

Original bias Slides

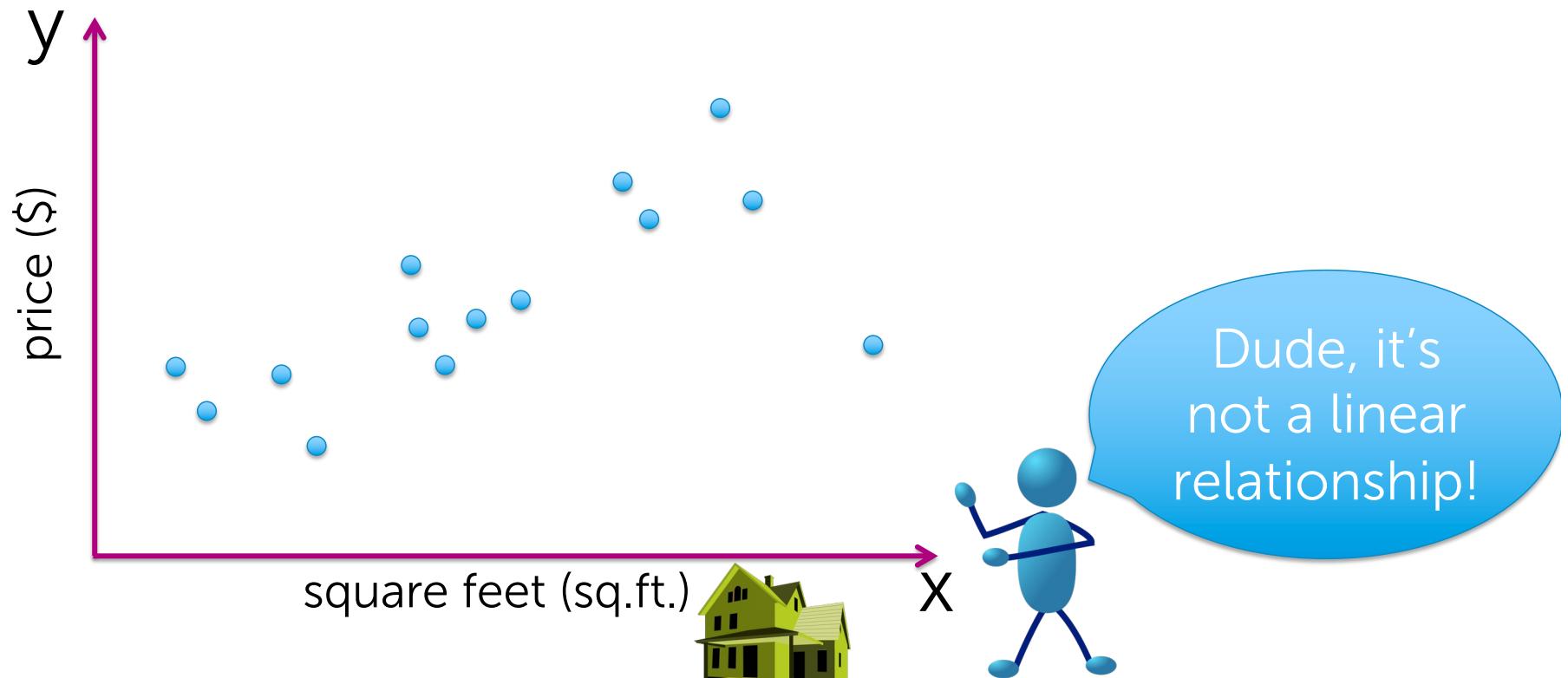
Bias-Variance Tradeoff



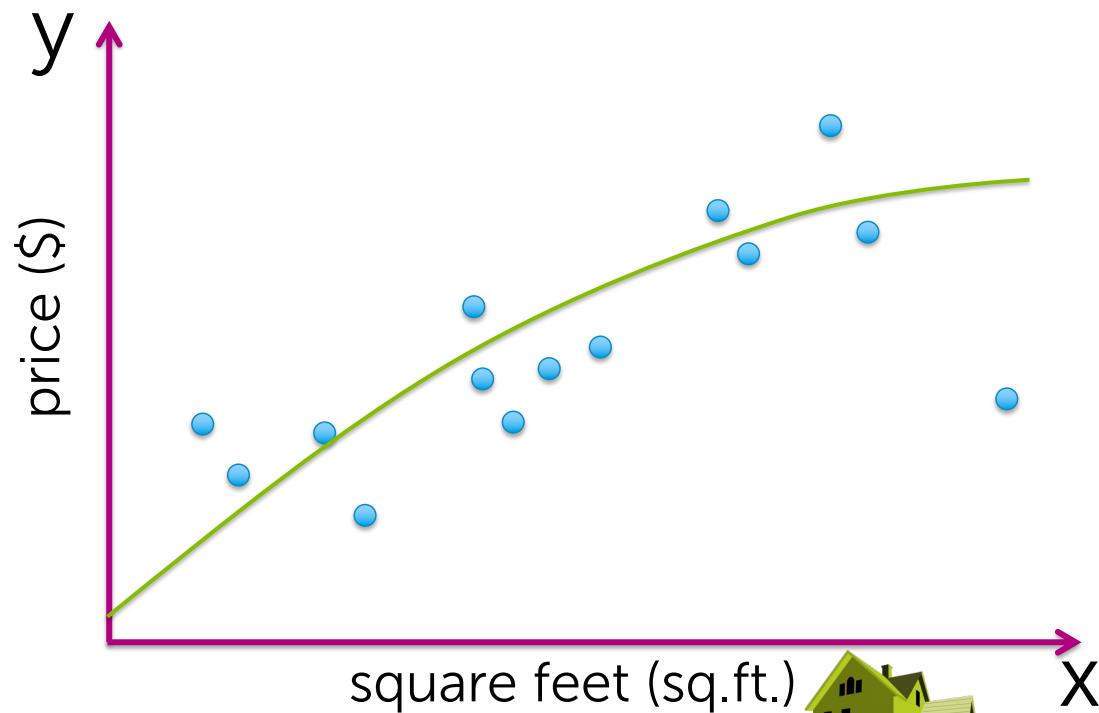
CS229: Machine Learning
Carlos Guestrin
Stanford University
Slides include content developed by and co-developed with Emily Fox



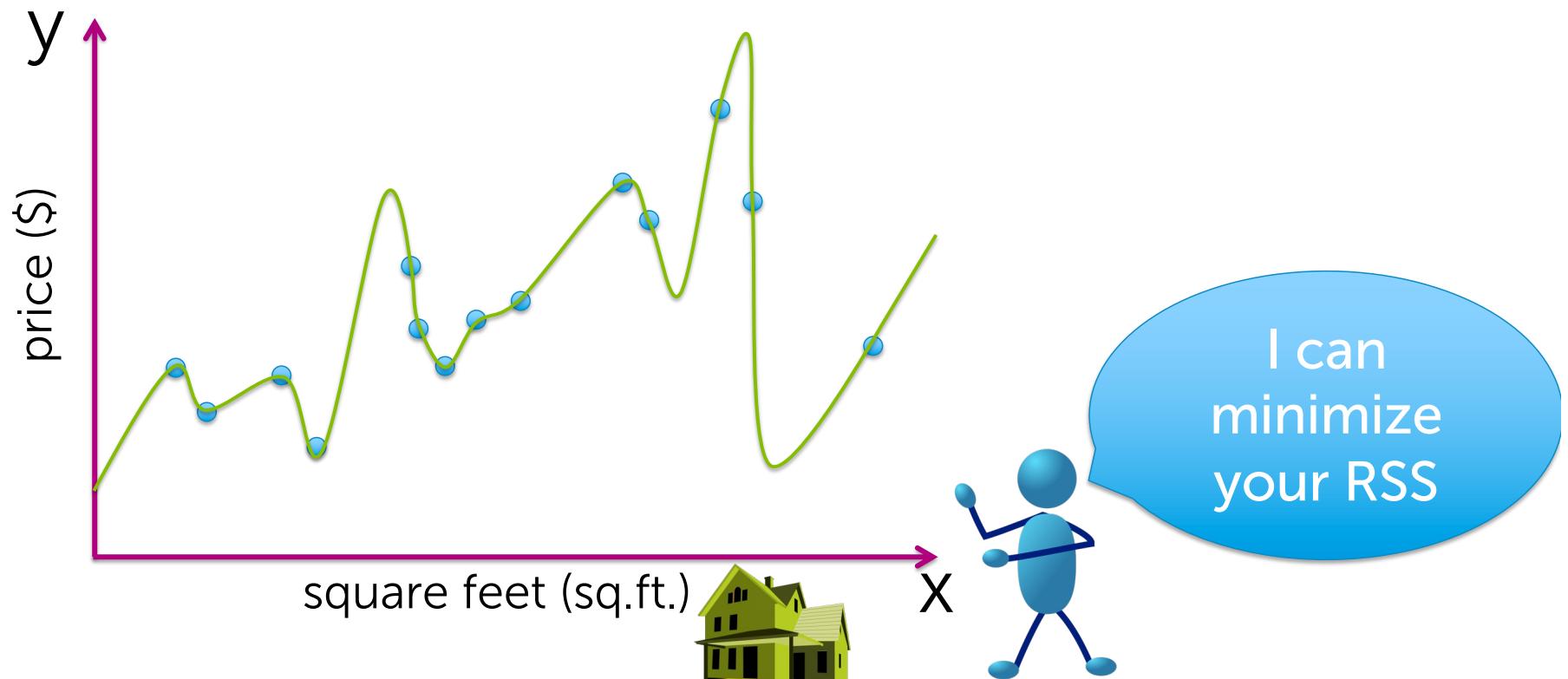
Fit data with a line or ... ?



What about a quadratic function?



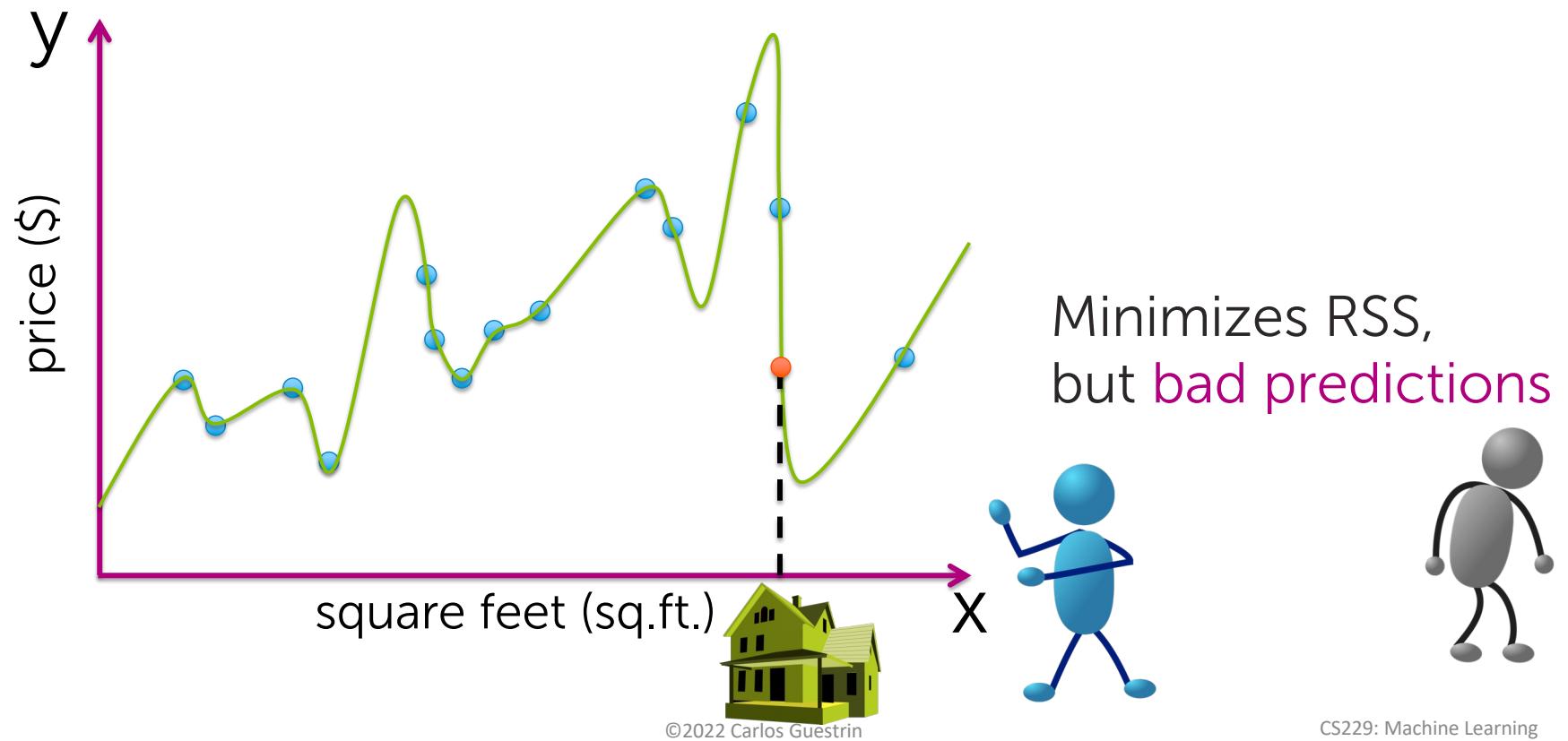
Even higher order polynomial



Do you believe this fit?



Do you believe this fit?



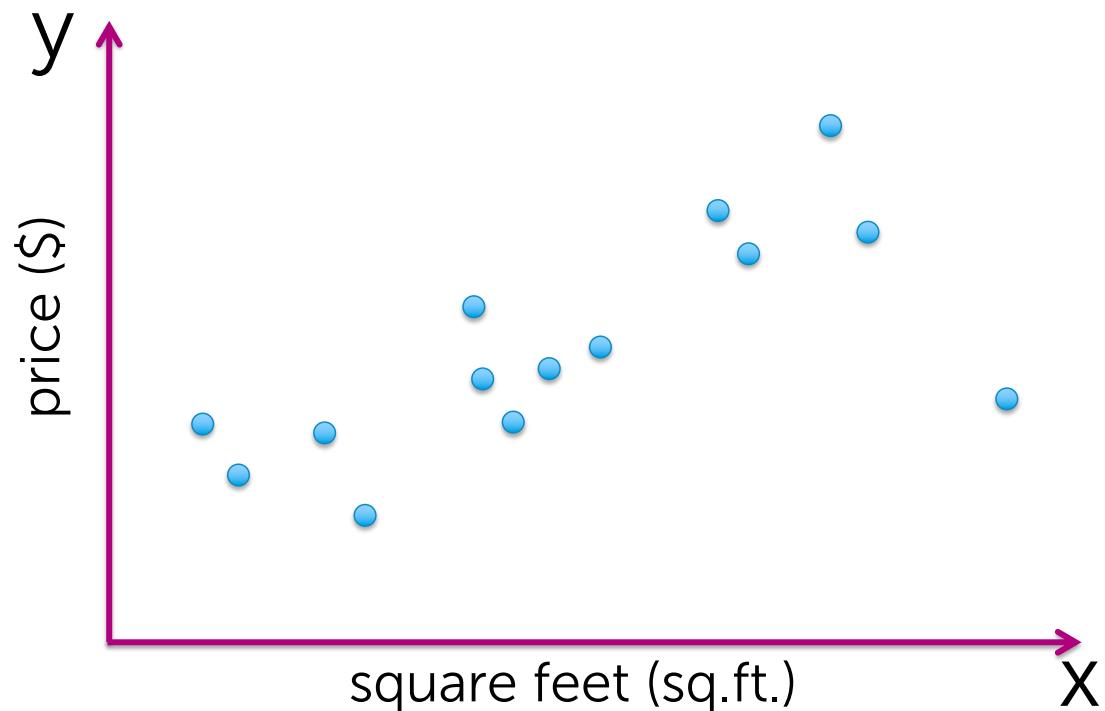
"Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful." George Box, 1987.

Assessing the loss

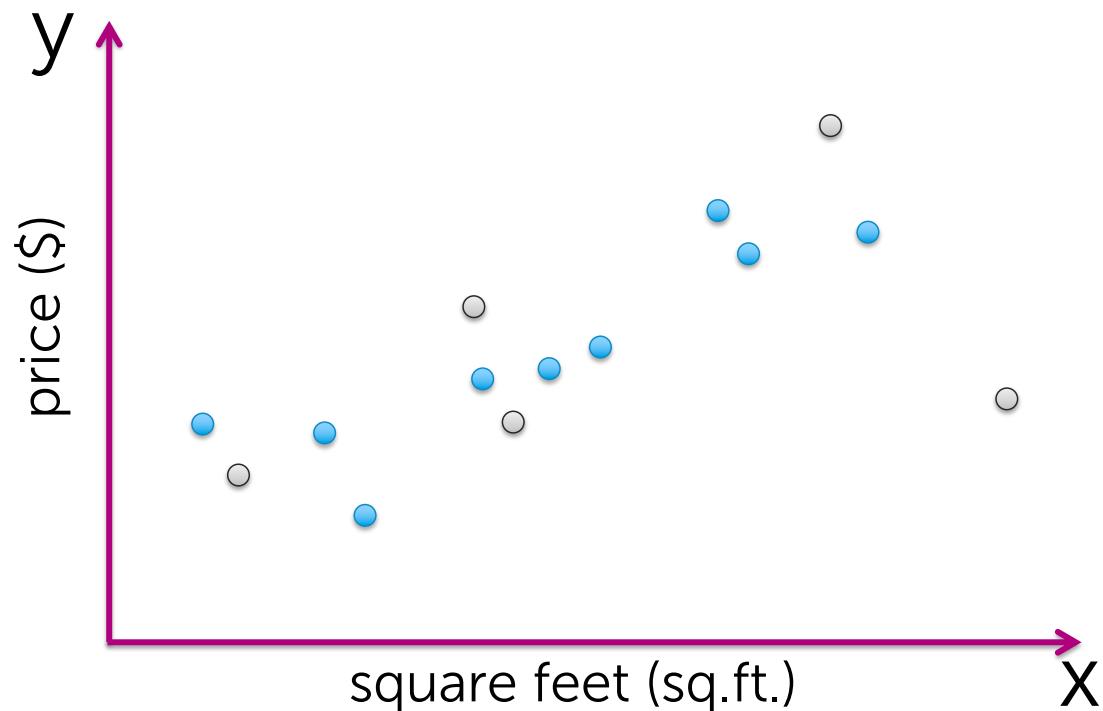
Assessing the loss

Part 1: Training error

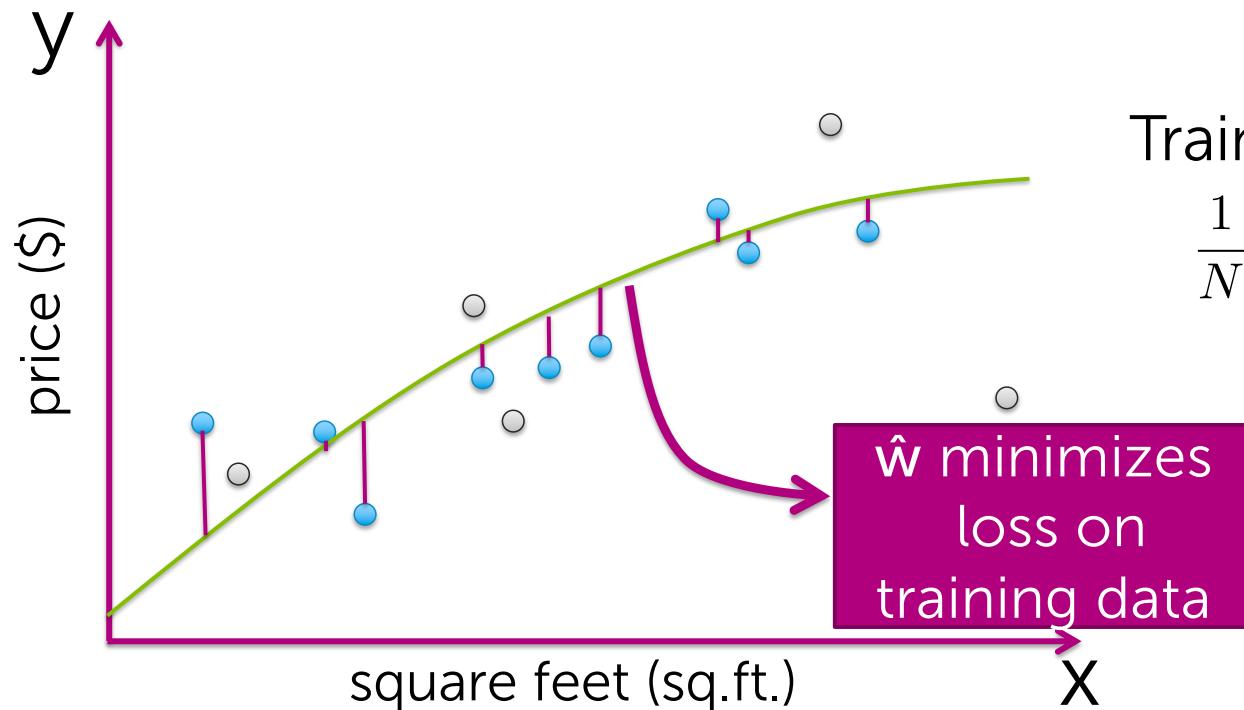
Define training data



Define training data



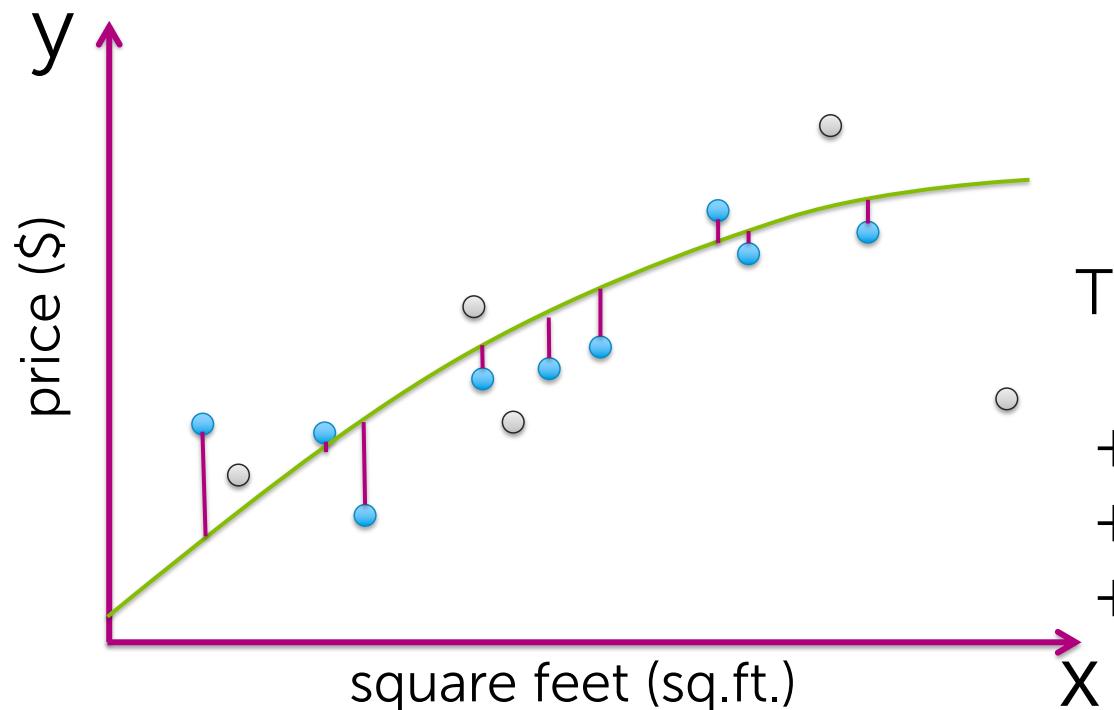
Example: Fit quadratic to minimize RSS



$$\text{Training error } (\hat{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - f_{\hat{w}}(x_i))^2$$

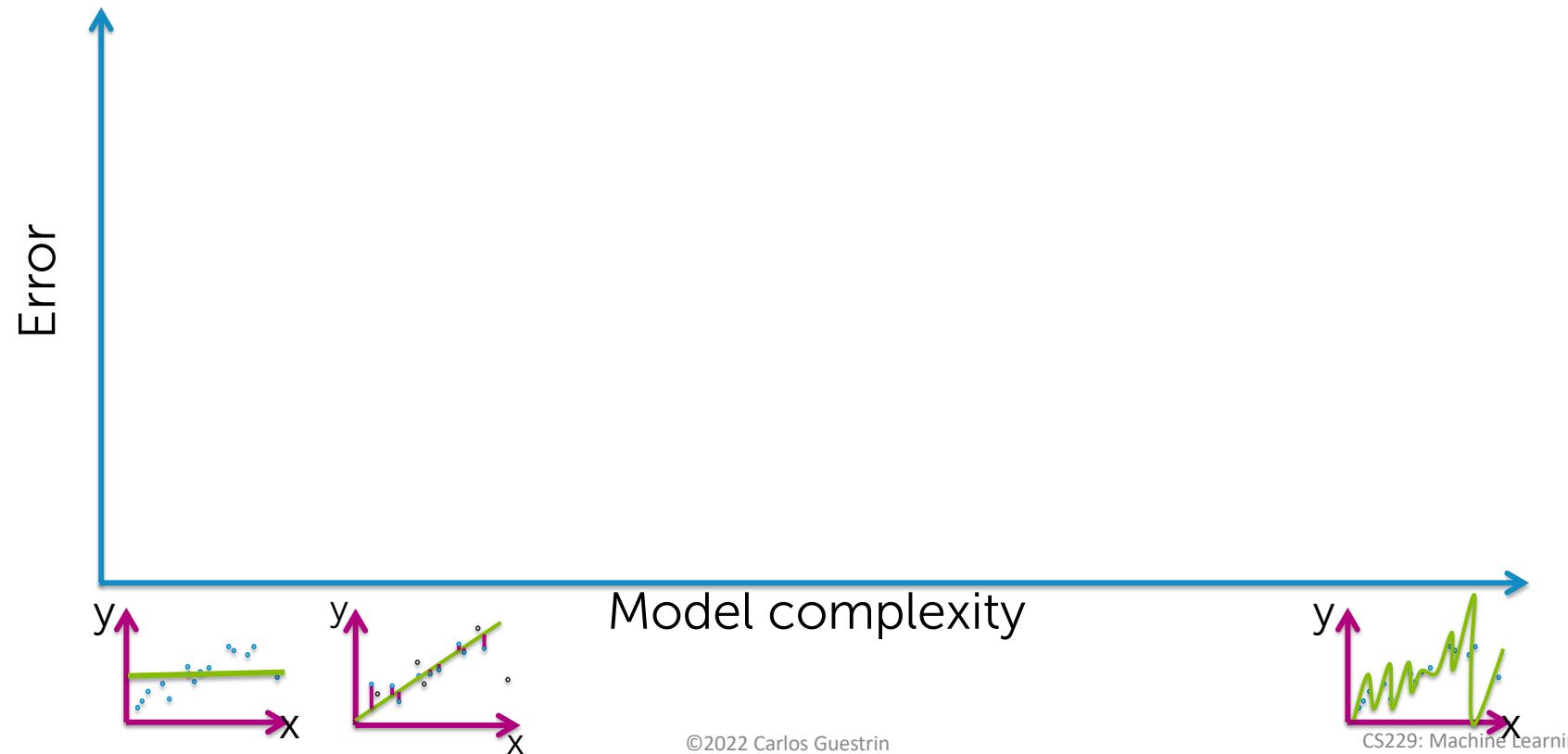
Example:

Use squared error loss $(y - f_{\hat{w}}(x))^2$



Training error (\hat{w}) = $1/N * [(\$_{train\ 1} - f_{\hat{w}}(\text{sq.ft.}_{train\ 1}))^2 + (\$_{train\ 2} - f_{\hat{w}}(\text{sq.ft.}_{train\ 2}))^2 + (\$_{train\ 3} - f_{\hat{w}}(\text{sq.ft.}_{train\ 3}))^2 + \dots \text{ include all training houses}]$

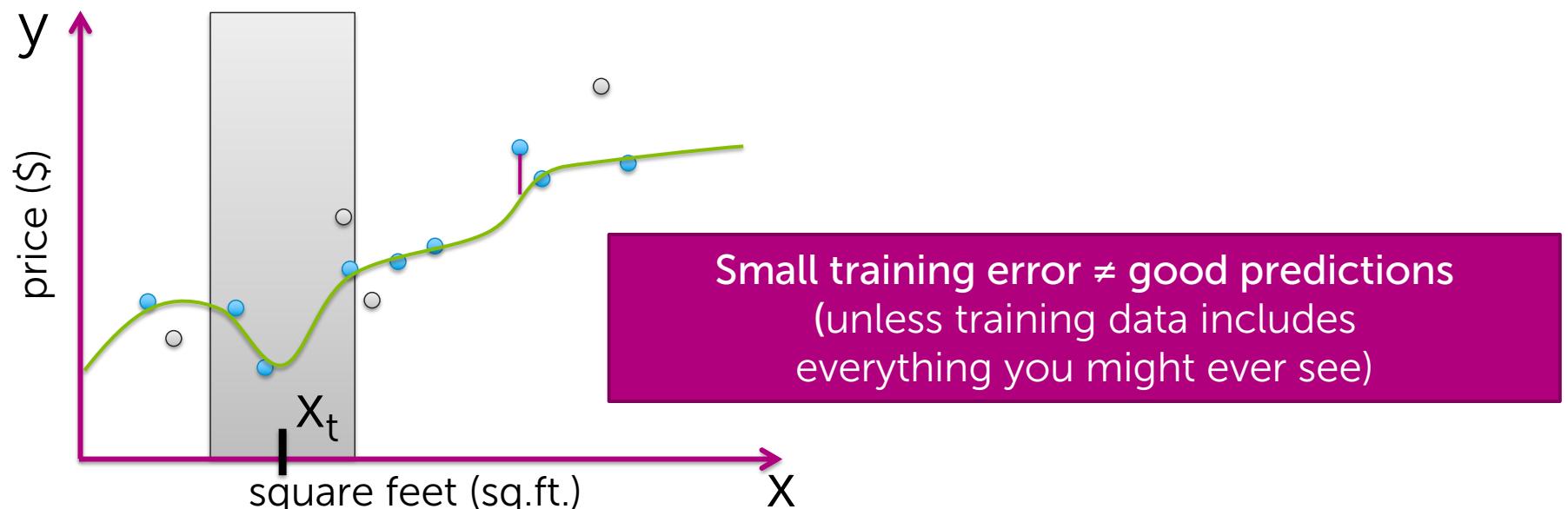
Training error vs. model complexity



Is training error a good measure of predictive performance?

Issue:

Training error is overly optimistic... \hat{w} was fit to training data

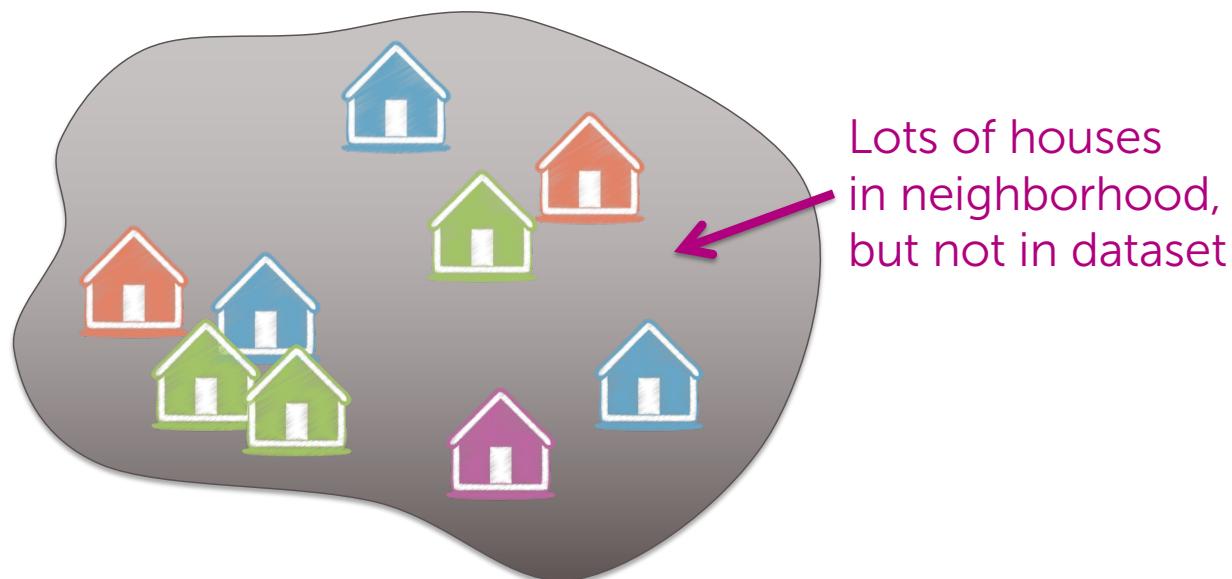


Assessing the loss

Part 2: Generalization (true) error

Generalization error

Really want estimate of loss over all possible (, ) pairs



Generalization error definition

Really want estimate of loss over all possible ( , \$) pairs

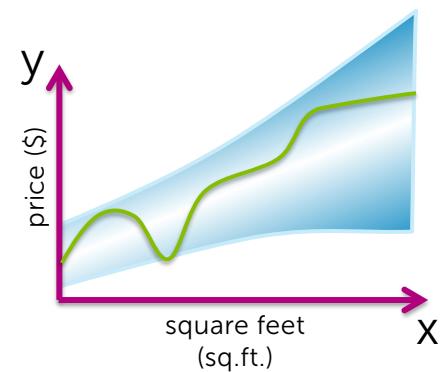
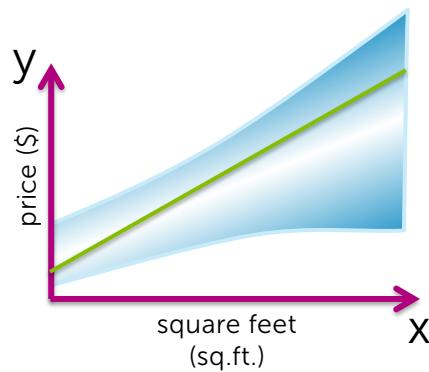
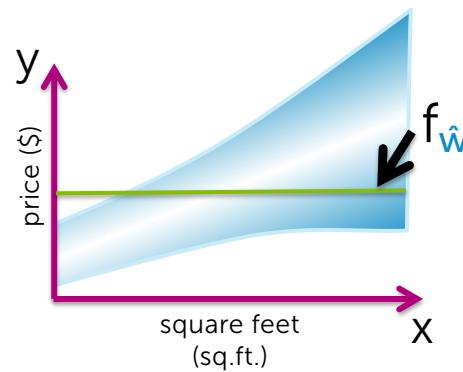
Formally:

average over all possible
(x, y) pairs weighted by
how likely each is

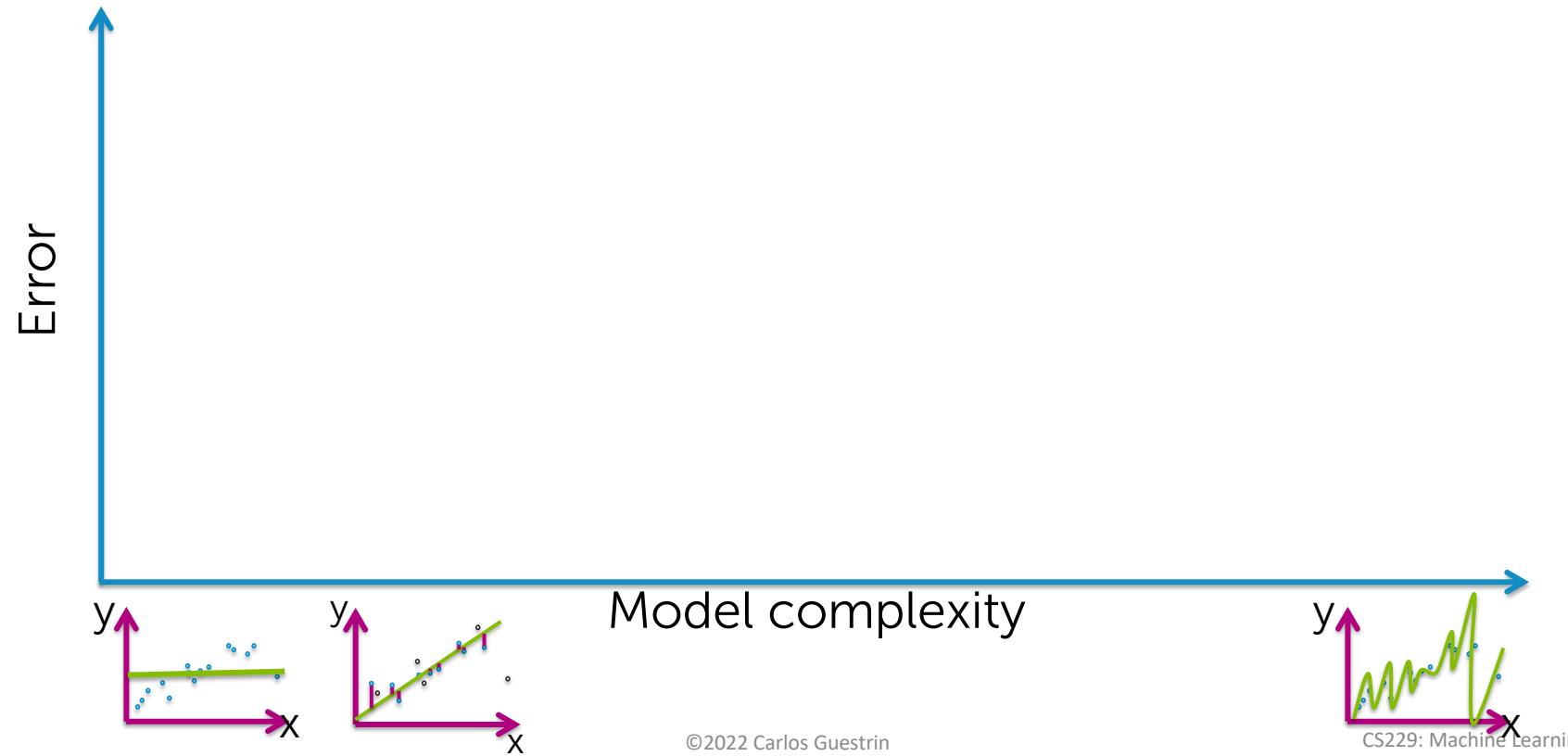
$$\text{generalization error} = E_{x,y}[L(y, f_{\hat{w}}(x))]$$

fit using training data

Generalization error vs. model complexity



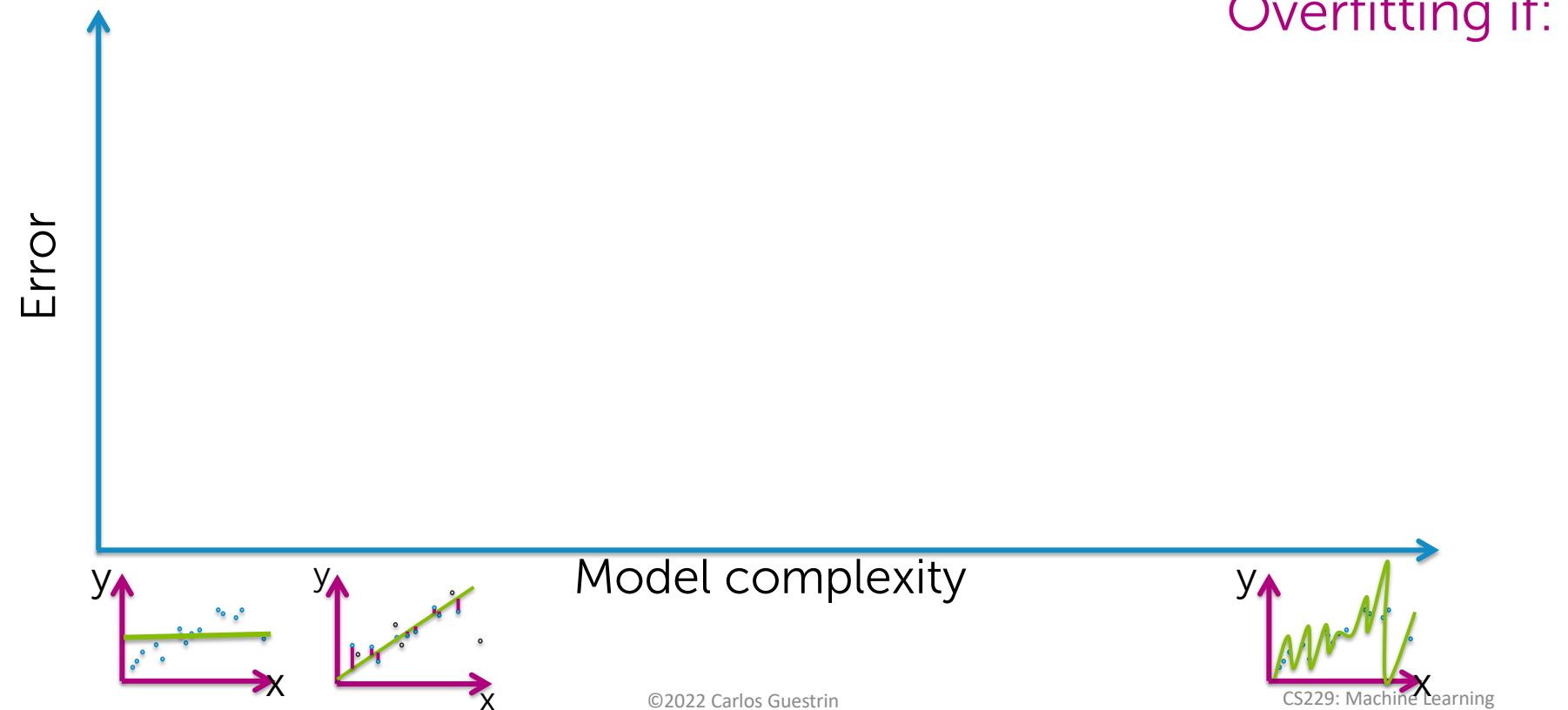
True error vs. model complexity



Assessing the loss

Part 3: Test error

Training, true, test error vs. model complexity



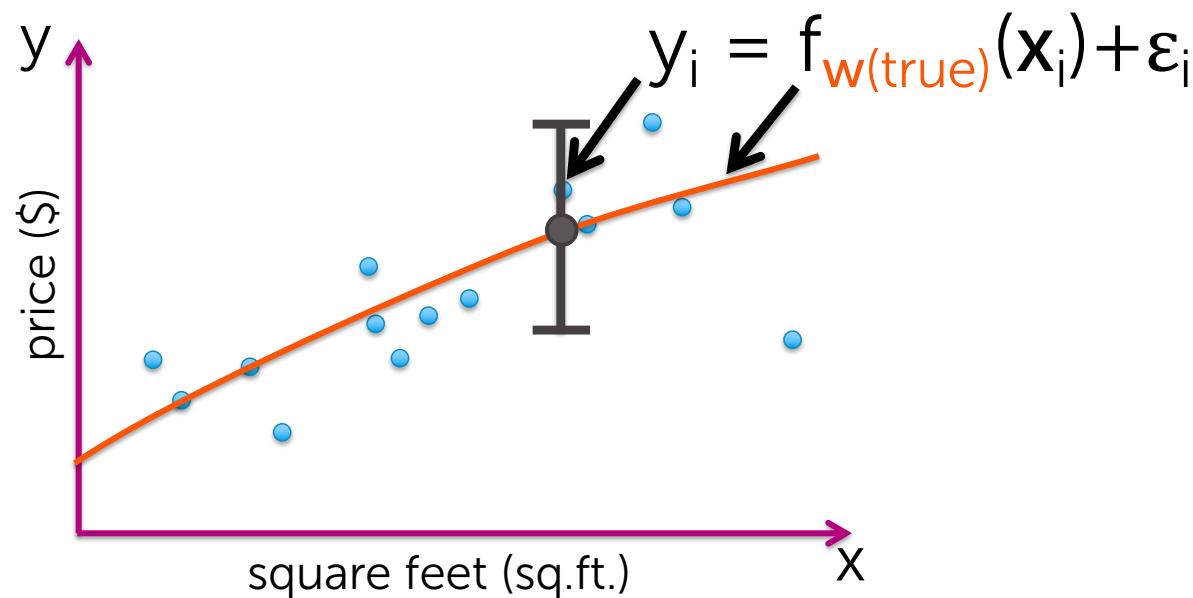
3 sources of error + the bias-variance tradeoff

3 sources of error

In forming predictions, there are 3 sources of error:

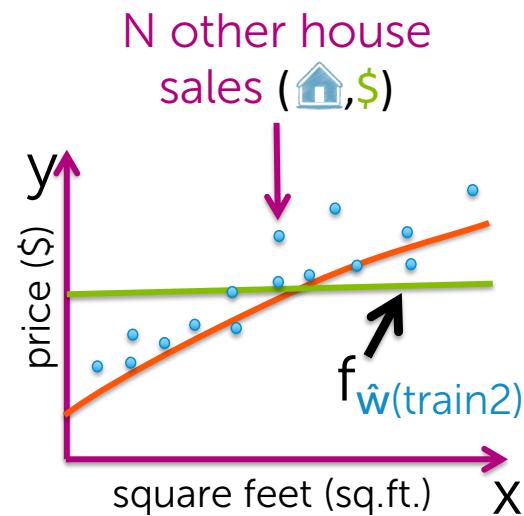
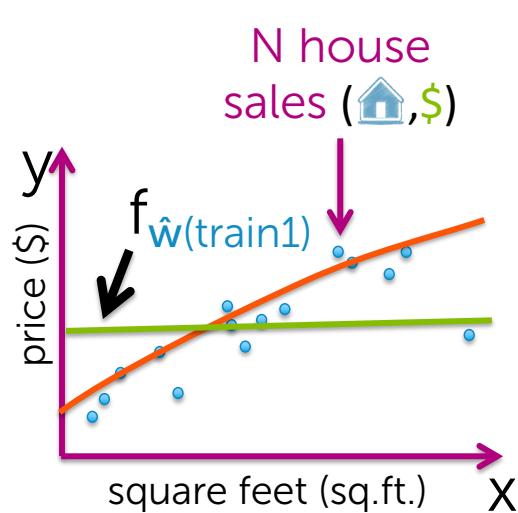
1. Noise
2. Bias
3. Variance

Data inherently **noisy**



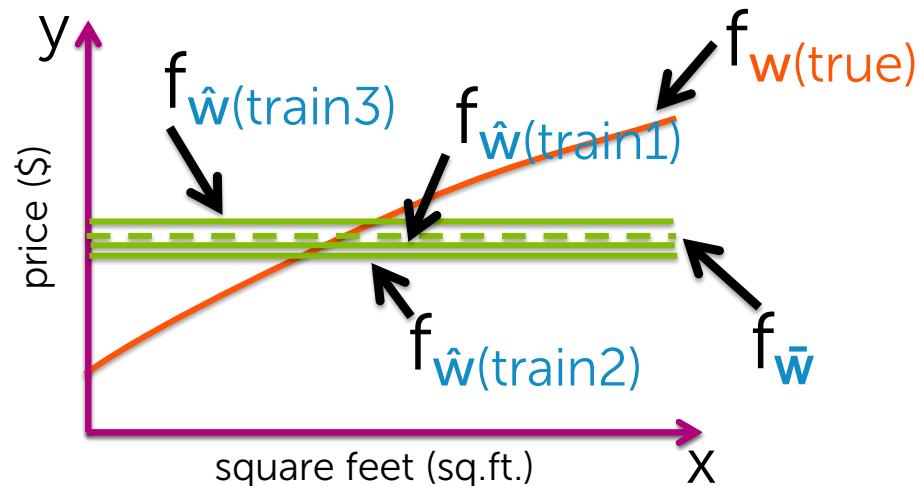
Bias contribution

Suppose we fit a constant function



Bias contribution

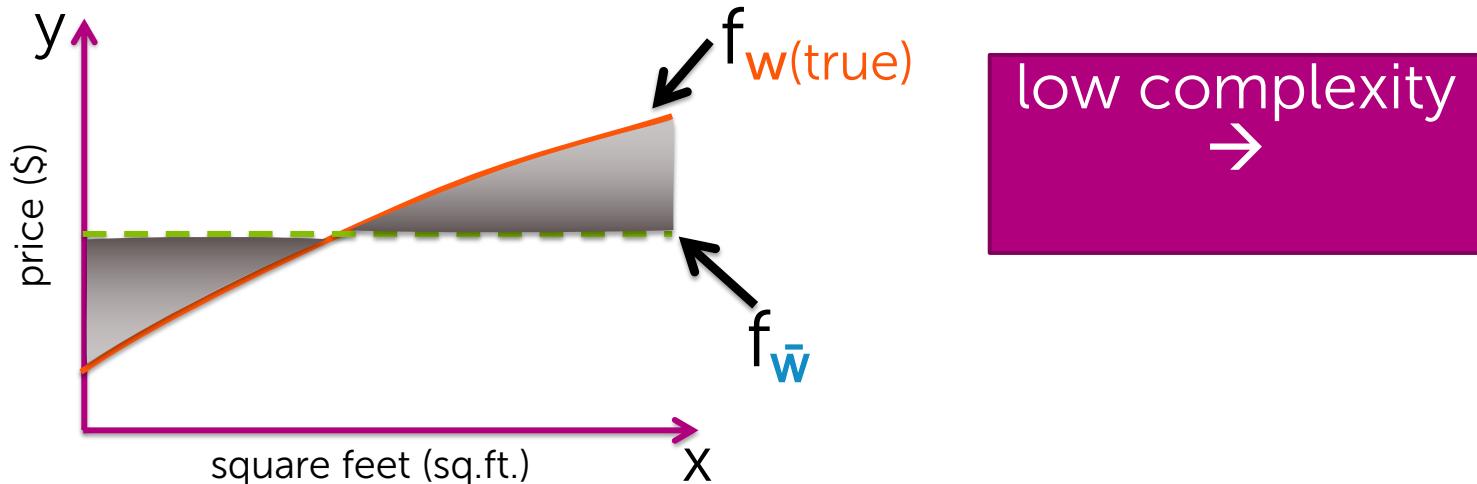
Over all possible size N training sets,
what do I expect my fit to be?



Bias contribution

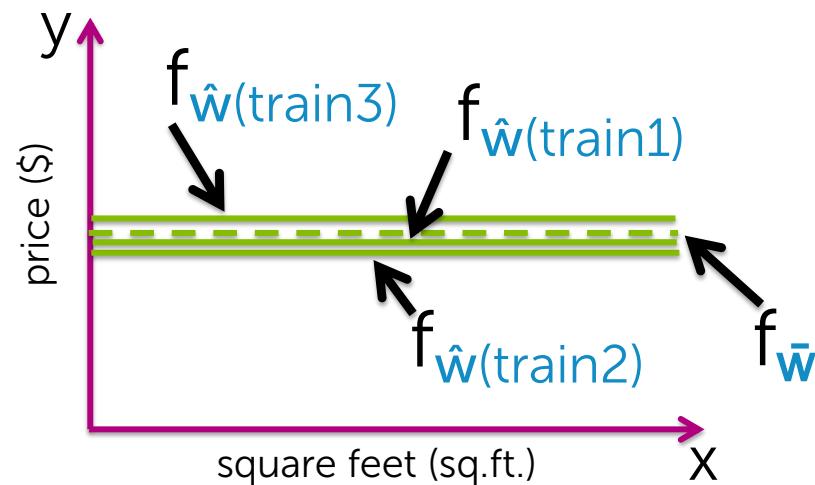
$$\text{Bias}(x) = f_{w(\text{true})}(x) - f_{\bar{w}}(x)$$

Is our approach flexible
enough to capture $f_{w(\text{true})}$?
If not, error in predictions.



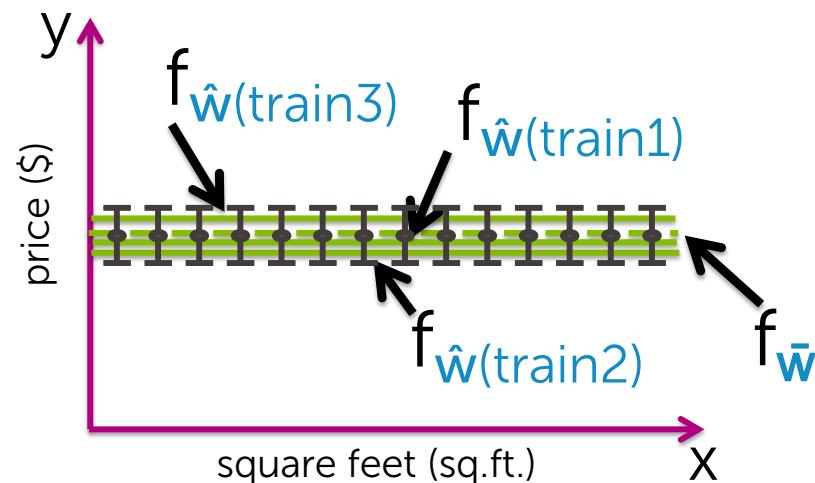
Variance contribution

How much do specific fits vary from the expected fit?



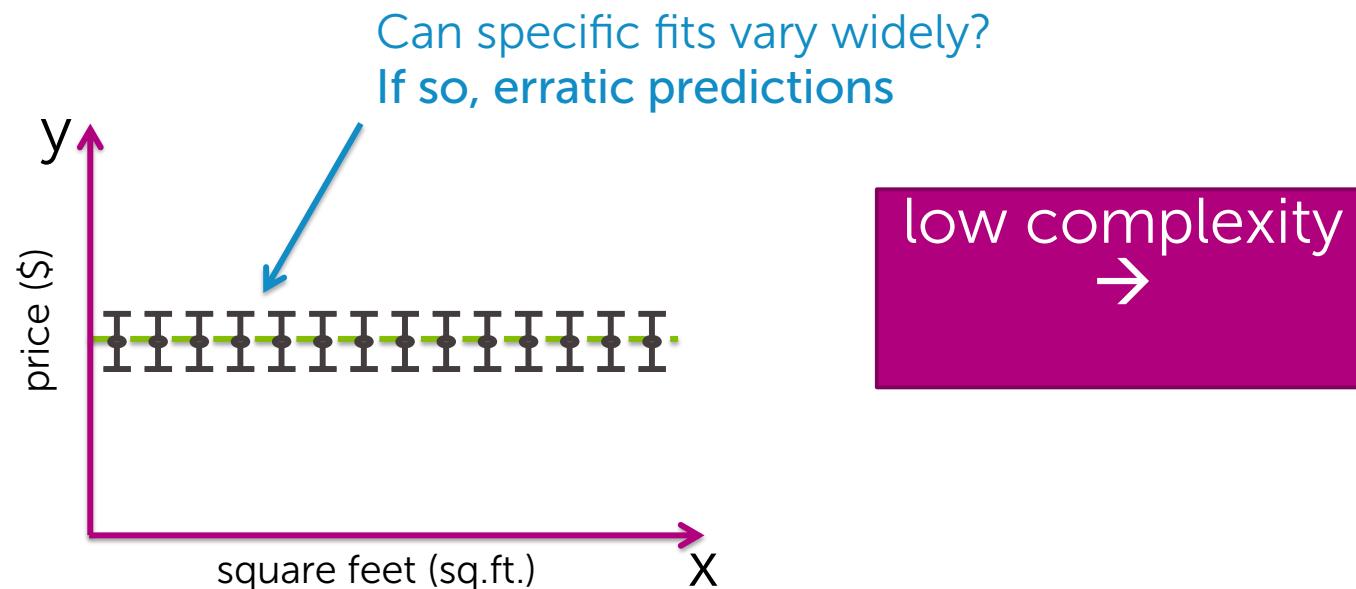
Variance contribution

How much do specific fits vary from the expected fit?



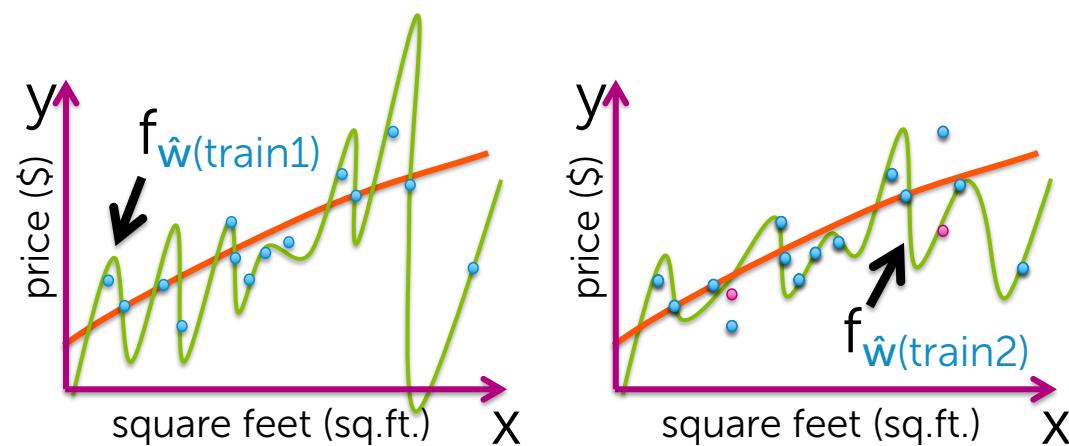
Variance contribution

How much do specific fits vary from the expected fit?



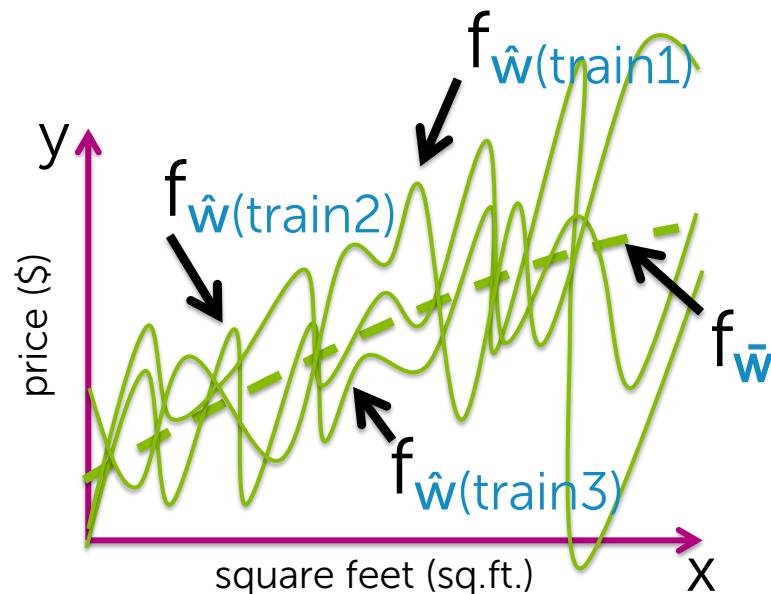
Variance of high-complexity models

Assume we fit a high-order polynomial

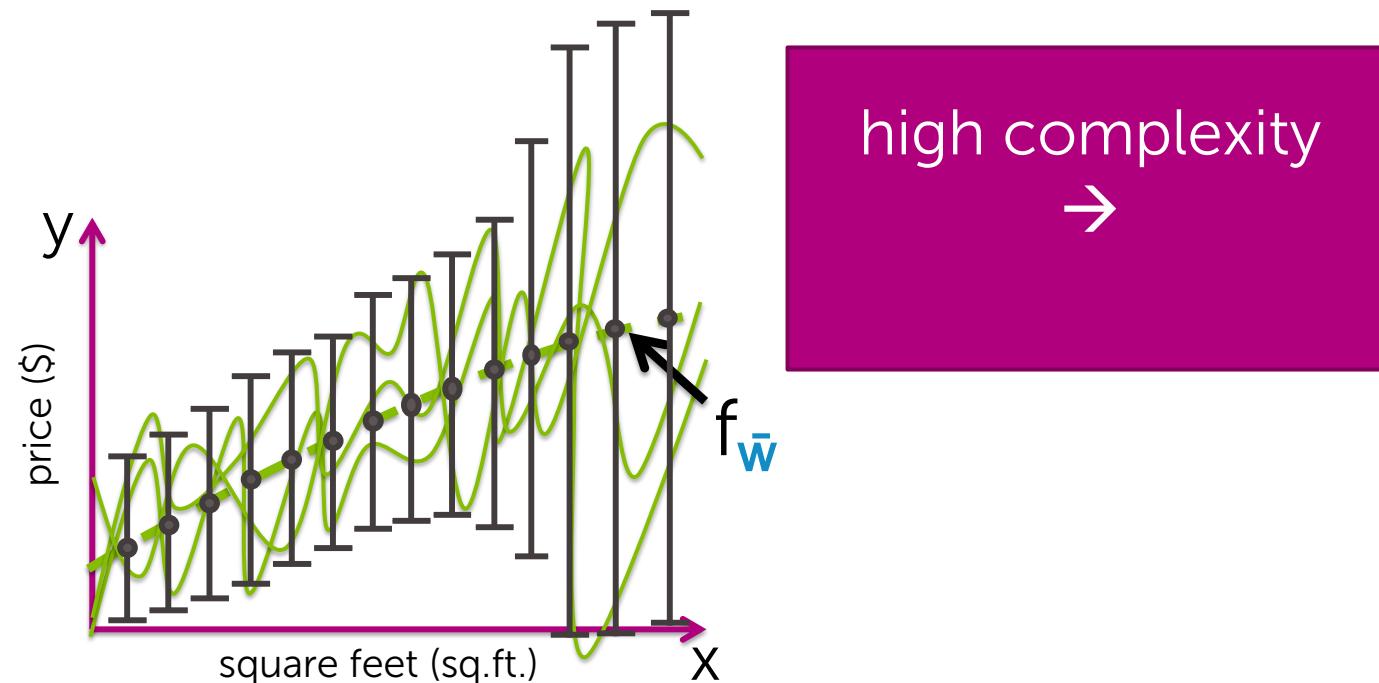


Variance of high-complexity models

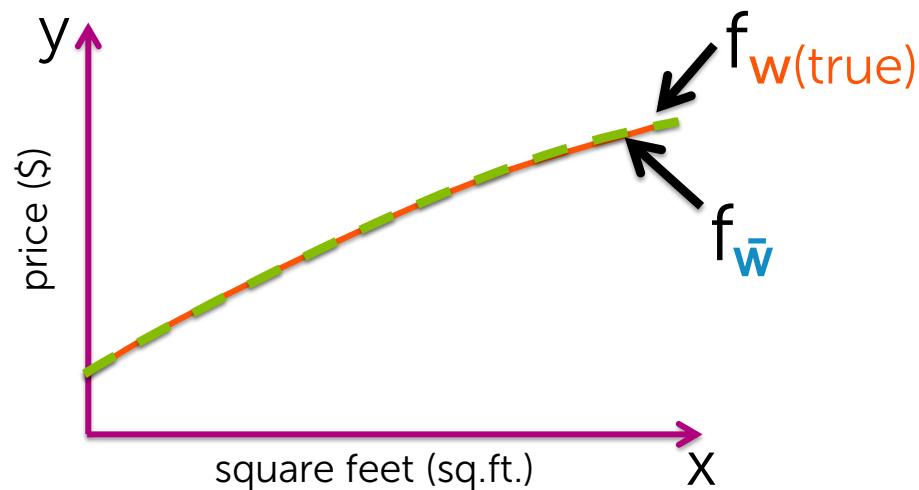
Suppose we fit a high-order polynomial



Variance of high-complexity models



Bias of high-complexity models



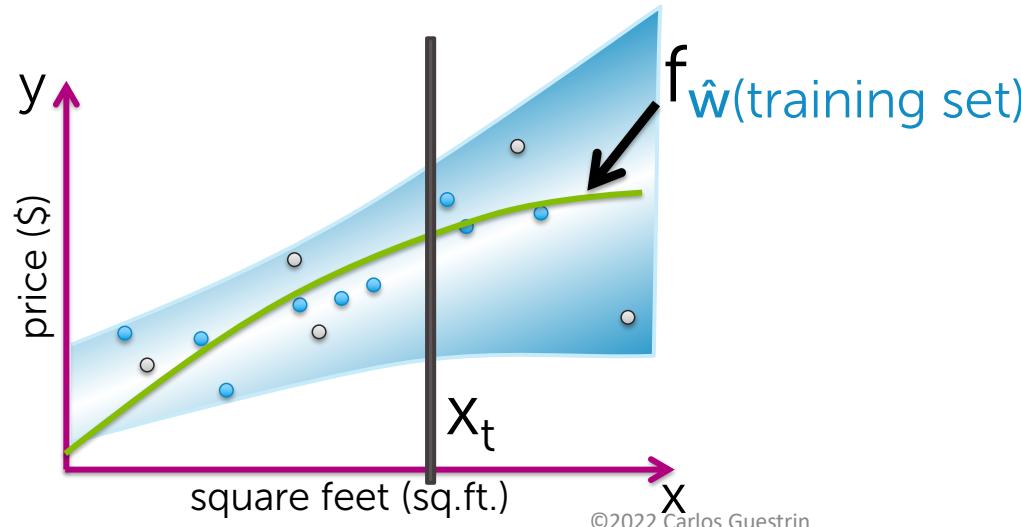
high complexity



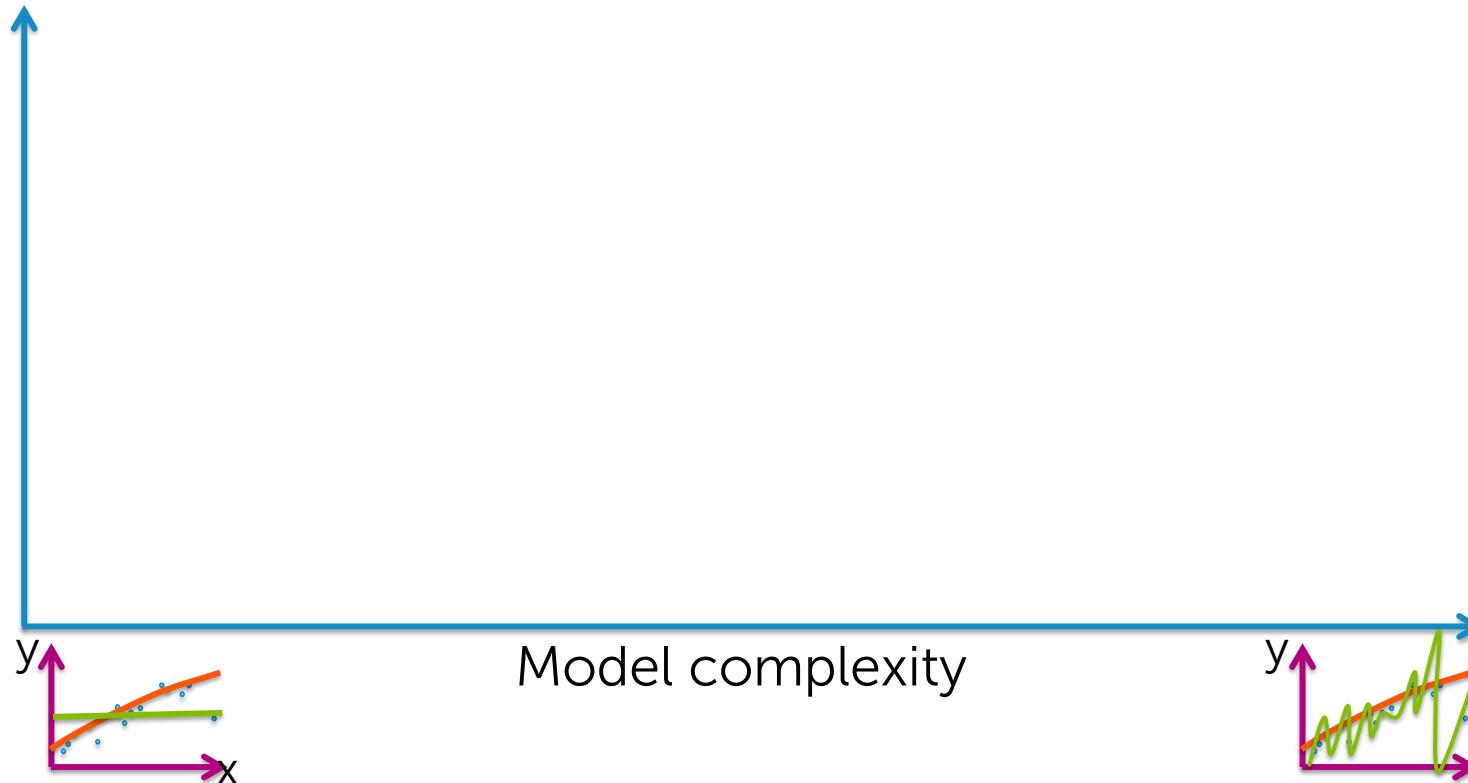
Sum of 3 sources of error

Average squared error at x_t

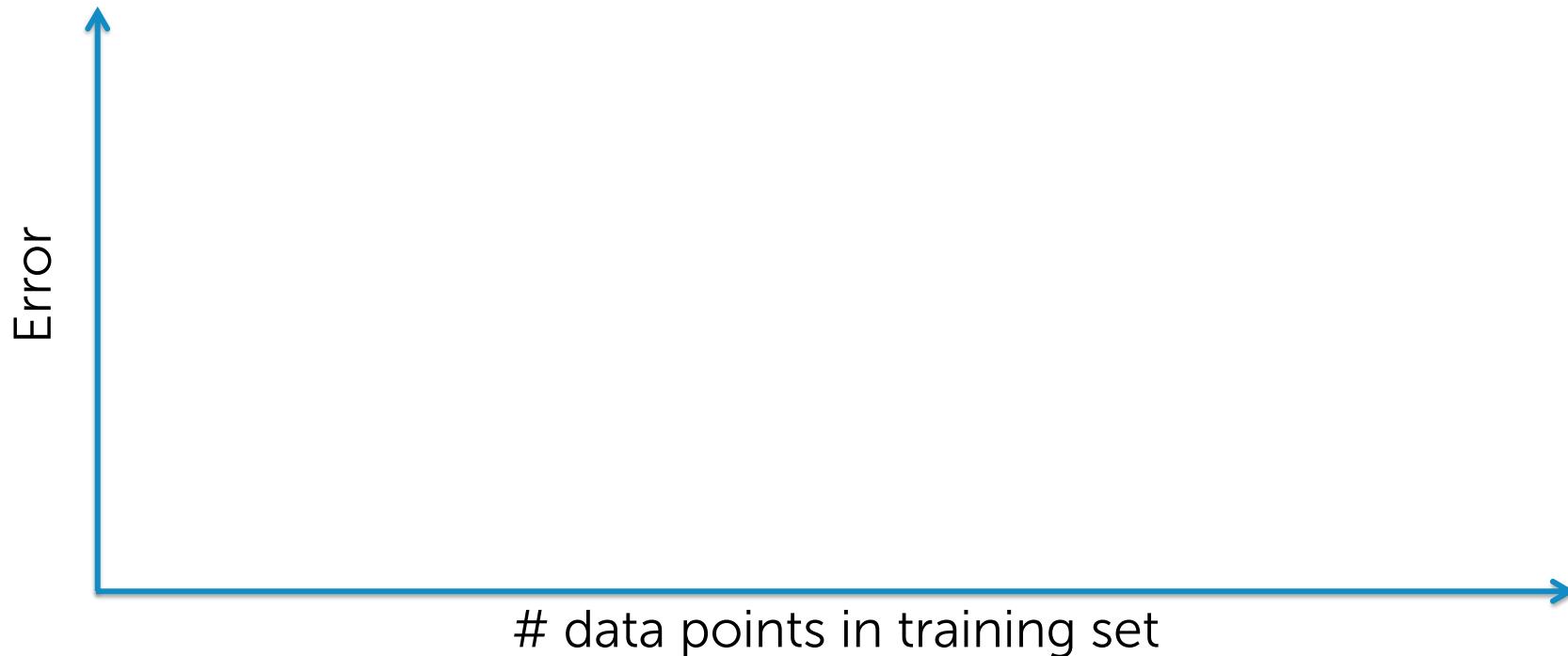
$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



Bias-variance tradeoff



Error vs. amount of data



Why 3 sources of error? A formal derivation

Deriving expected prediction error

Expected prediction error

$$\begin{aligned} &= E_{\text{train}} [\text{generalization error of } \hat{\mathbf{w}}(\text{train})] \\ &= E_{\text{train}} [E_{x,y} [L(y, f_{\hat{\mathbf{w}}(\text{train})}(x))]] \end{aligned}$$

1. Look at specific x_t
2. Consider $L(y, f_{\hat{\mathbf{w}}}(x)) = (y - f_{\hat{\mathbf{w}}}(x))^2$

Expected prediction error at x_t

$$= E_{\text{train}, y_t} [(y_t - f_{\hat{\mathbf{w}}(\text{train})}(x_t))^2]$$

Simplifying Notation

- Expected prediction error at x_t
 $= E_{\text{train}, y_t} [(y_t - f_{\hat{w}(\text{train})}(x_t))^2]$
- Simple (and abusive 😊) notation:
 - $y_t \rightarrow y$
 - $f_{w(\text{true})}(x_t) \rightarrow f$
 - $f_{\hat{w}(\text{train})}(x_t) \rightarrow \hat{f}$
 - $E_{\text{train}}[f_{\hat{w}(\text{train})}(x_t)] = f_{\bar{w}}(x_t) \rightarrow \bar{f}$

Deriving expected prediction error

Expected prediction error at x_t

$$\begin{aligned} &= E_{\text{train}, y_t} [(y_t - f_{\hat{w}(\text{train})}(x_t))^2] = E_{\text{train}} [(y - \hat{f})^2] = \\ &= E_{\text{train}} [(y - f) + (f - \hat{f})]^2 \end{aligned}$$

Equating MSE with bias and variance

$$\begin{aligned}\text{MSE}[f_{\hat{w}(\text{train})}(x_t)] \\ &= E_{\text{train}}[(f - \hat{f})^2] \\ &= E_{\text{train}}[((f - \bar{f}) + (\bar{f} - \hat{f}))^2]\end{aligned}$$

Putting it all together

Expected prediction error at x_t

$$= \sigma^2 + \text{MSE}[f_{\hat{w}}(x_t)]$$

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



3 sources of error

Summary of bias-variance tradeoff

What you can do now...

- Contrast relationship between model complexity and train, true and test loss
- Compute training and test error given a loss function for different model complexities
- List and interpret the 3 sources of avg. prediction error
 - Irreducible error, bias, and variance

Annotated bias Slides

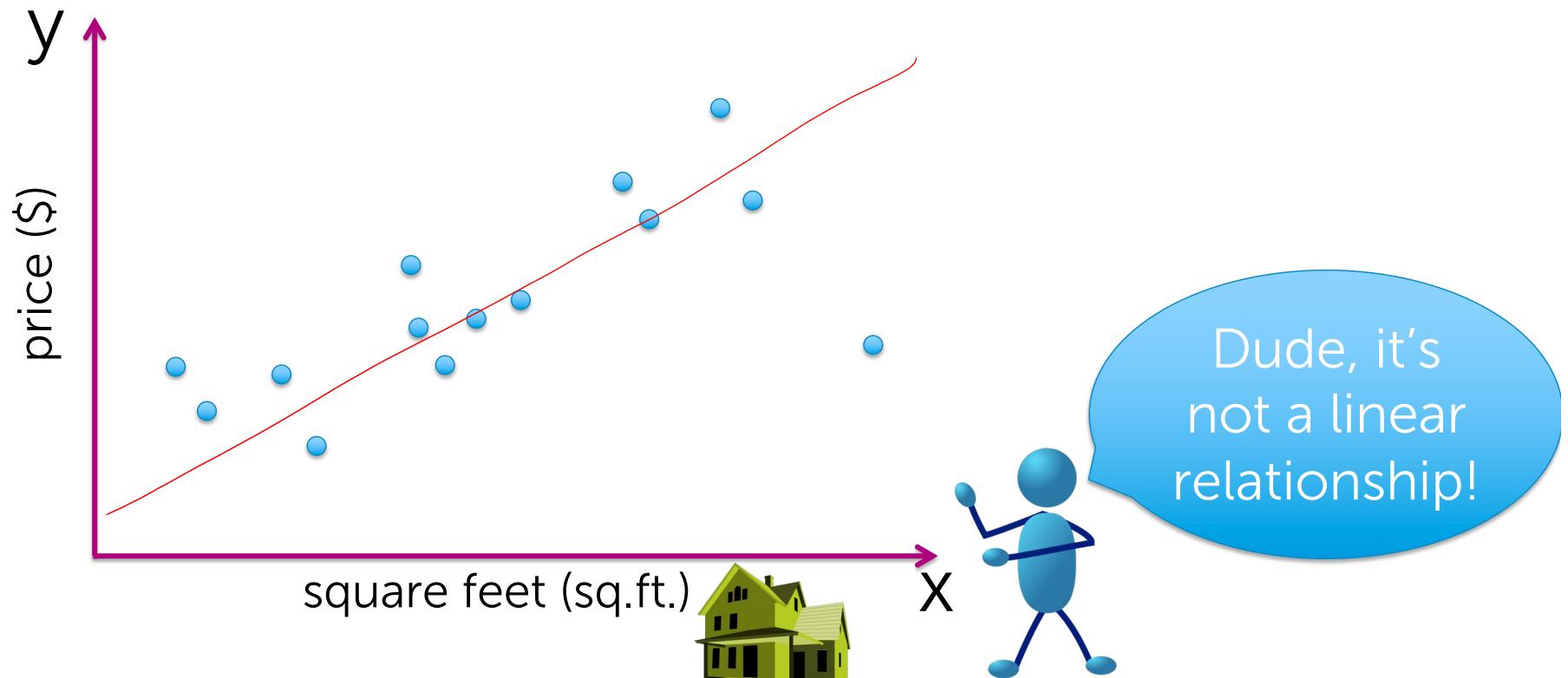
Bias-Variance Tradeoff



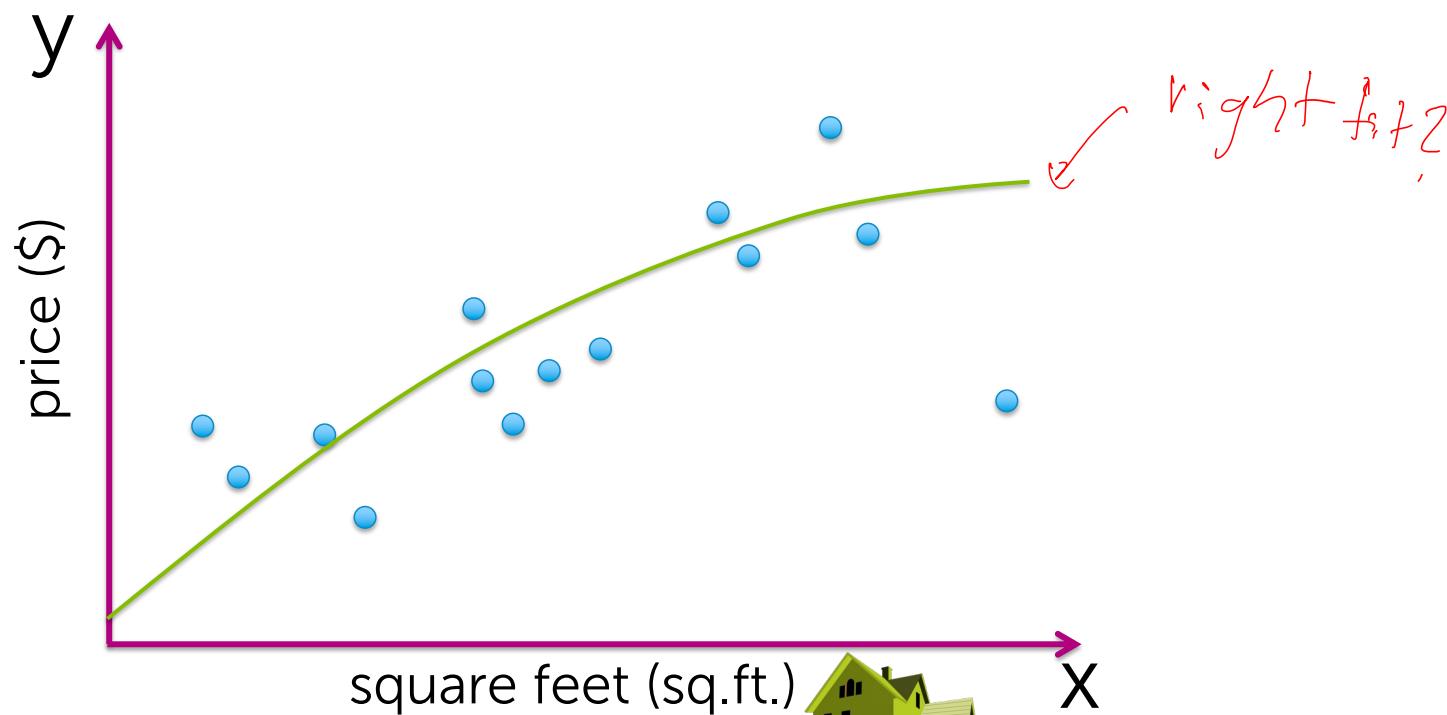
CS229: Machine Learning
Carlos Guestrin
Stanford University
Slides include content developed by and co-developed with Emily Fox



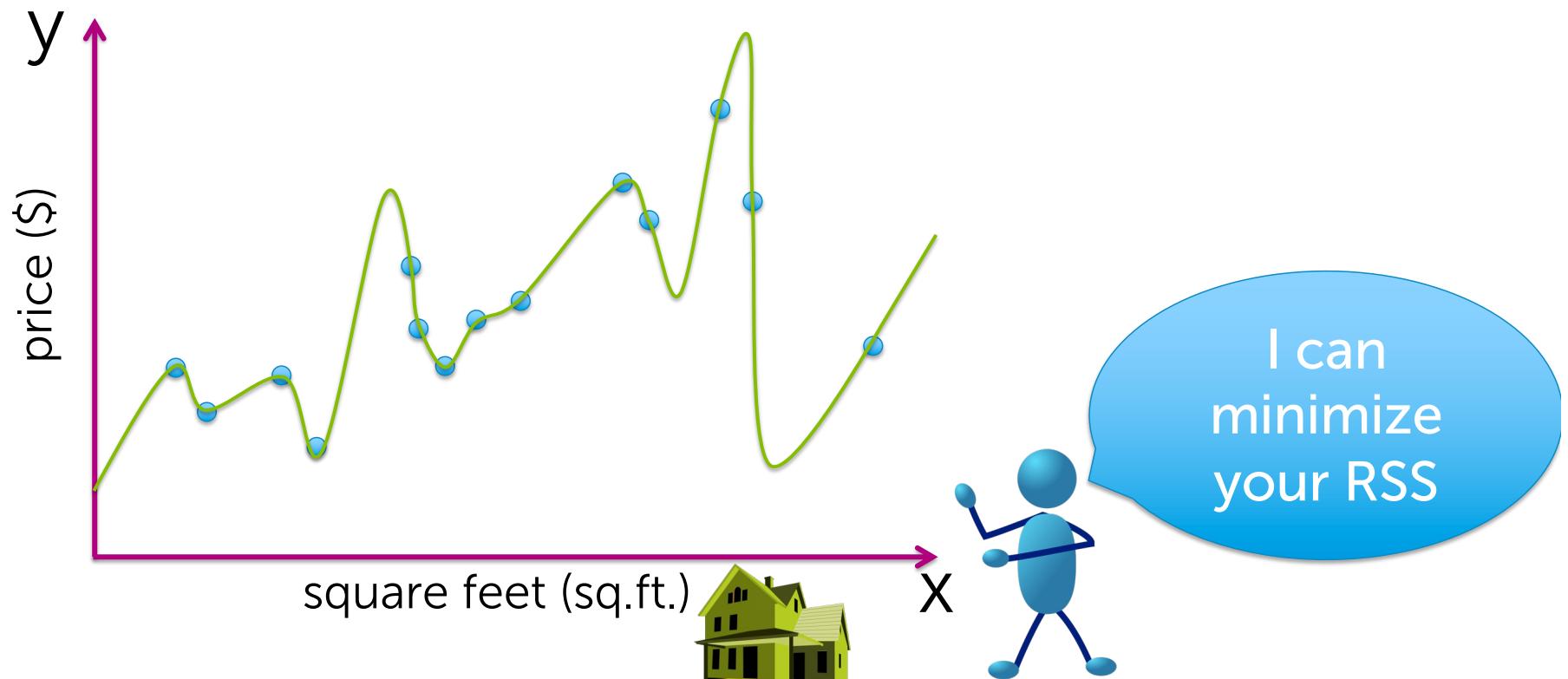
Fit data with a line or ... ?



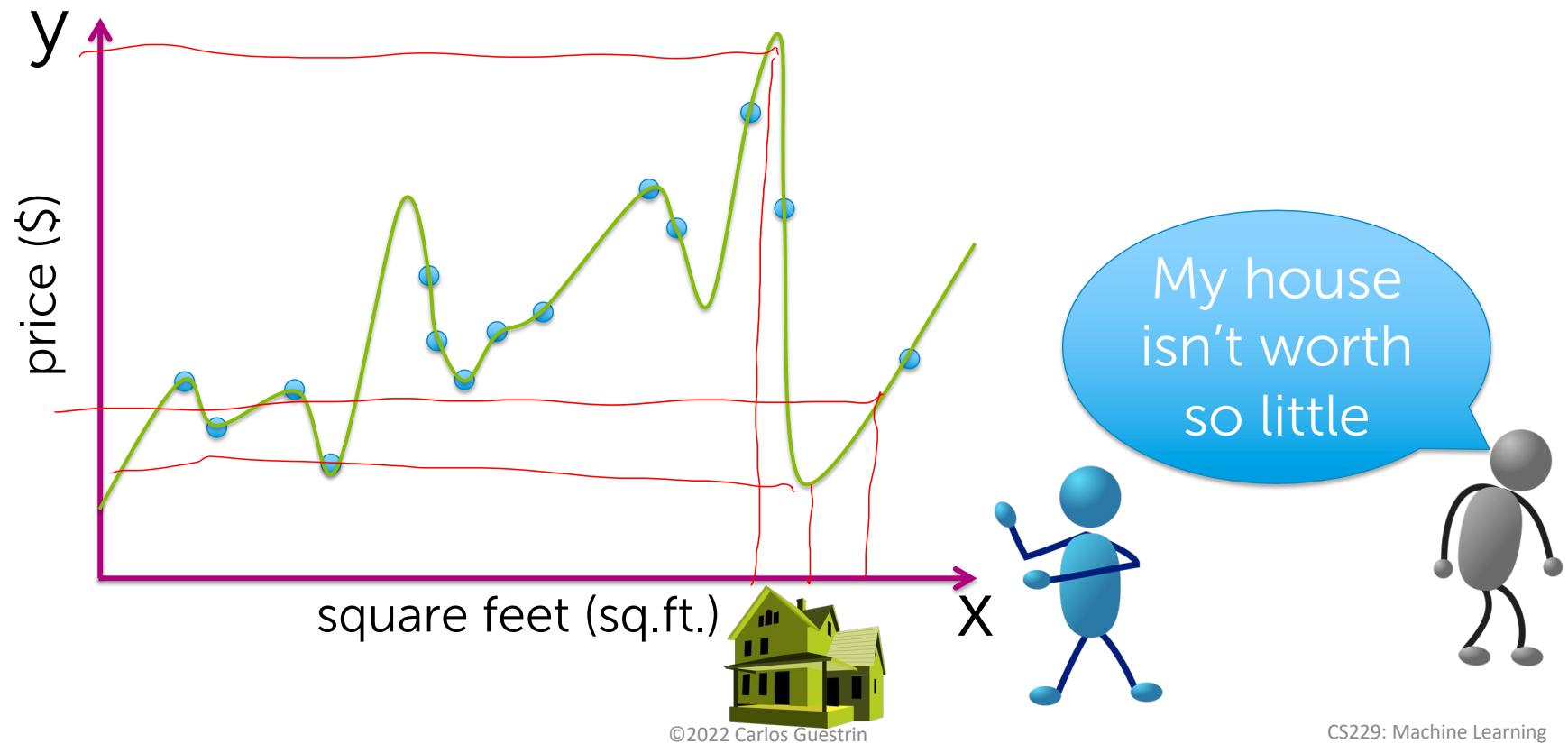
What about a quadratic function?



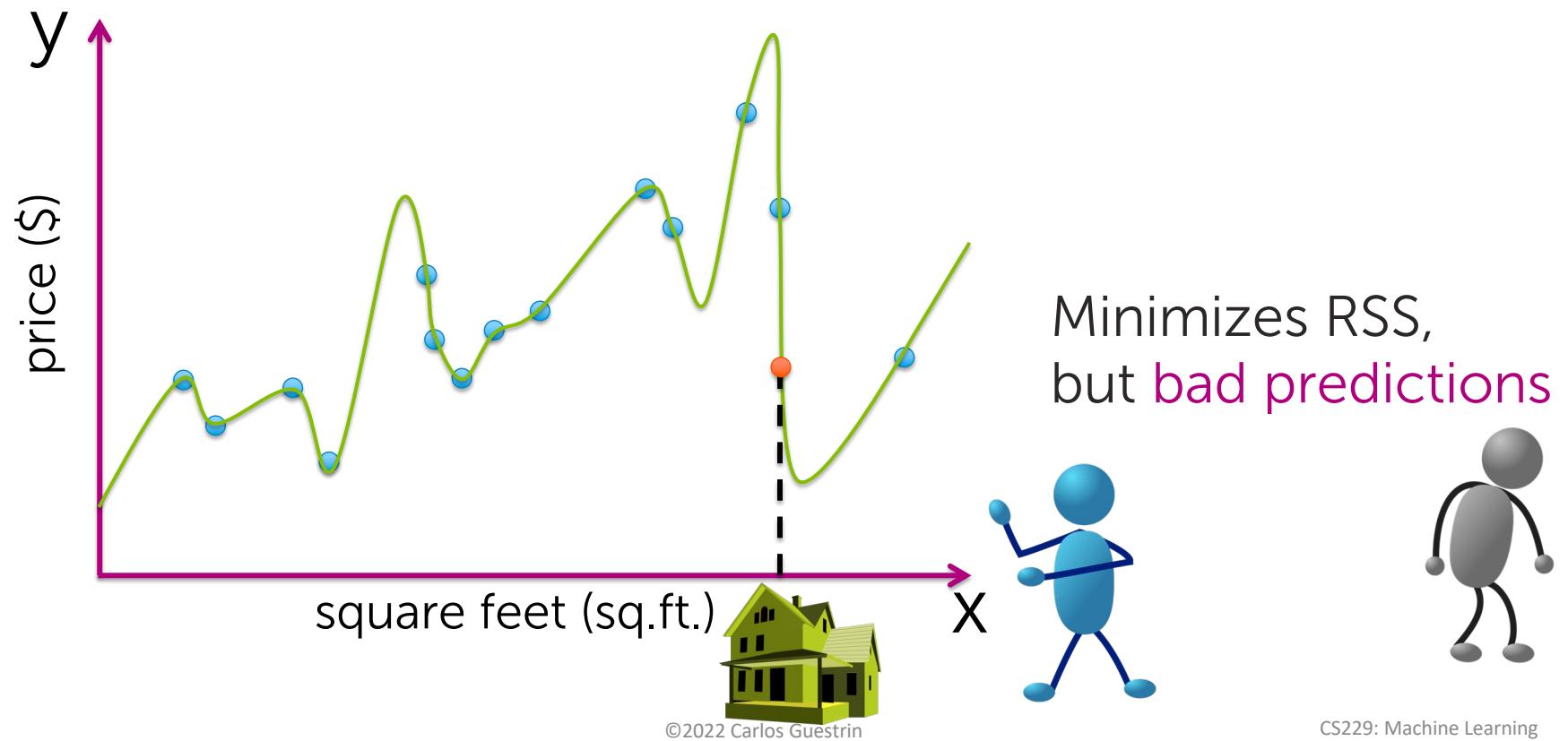
Even higher order polynomial



Do you believe this fit?



Do you believe this fit?



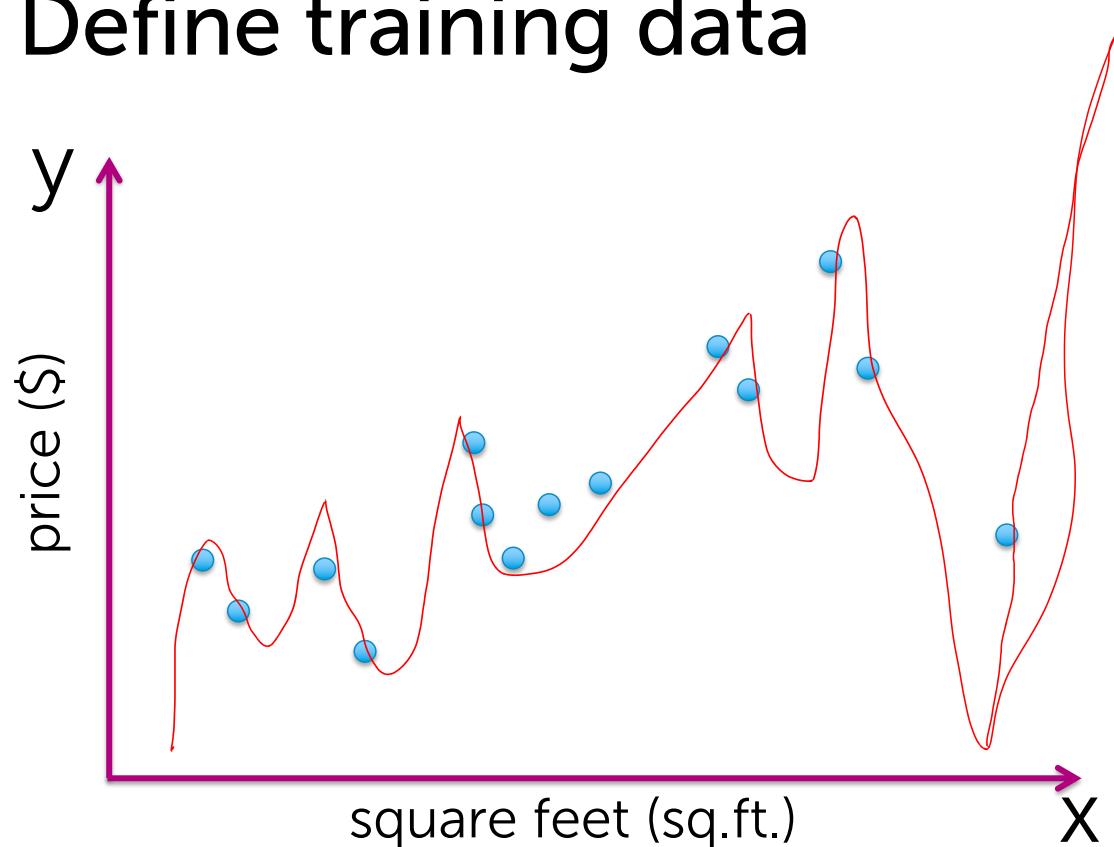
"Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful." George Box, 1987.

Assessing the loss

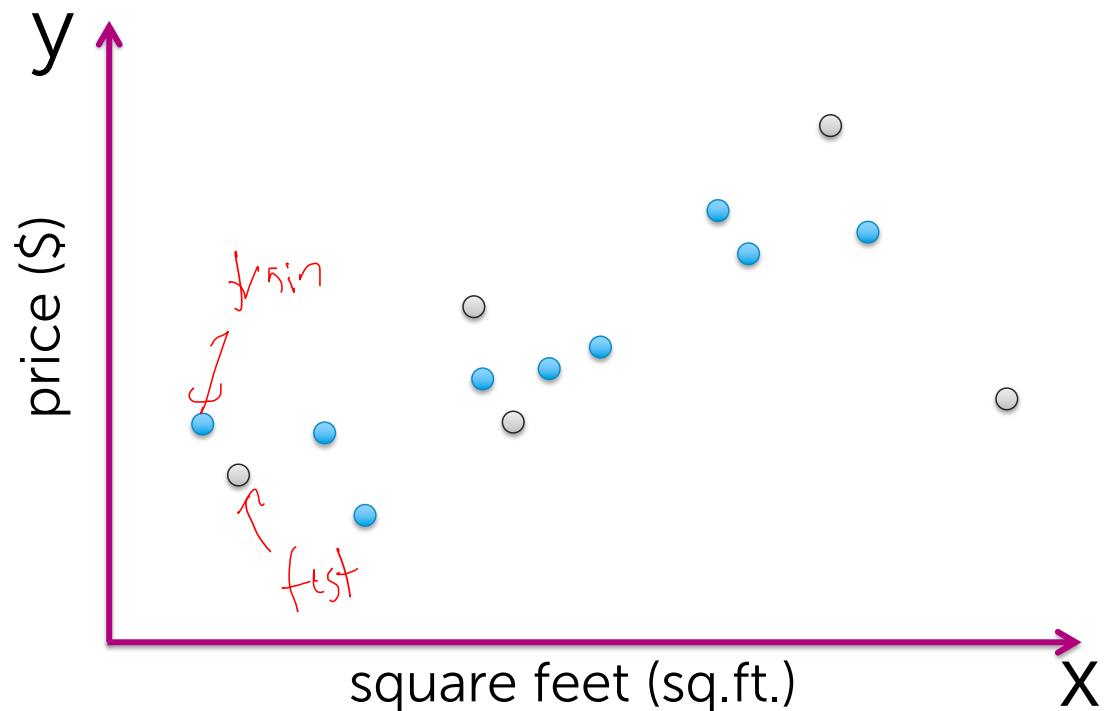
Assessing the loss

Part 1: Training error

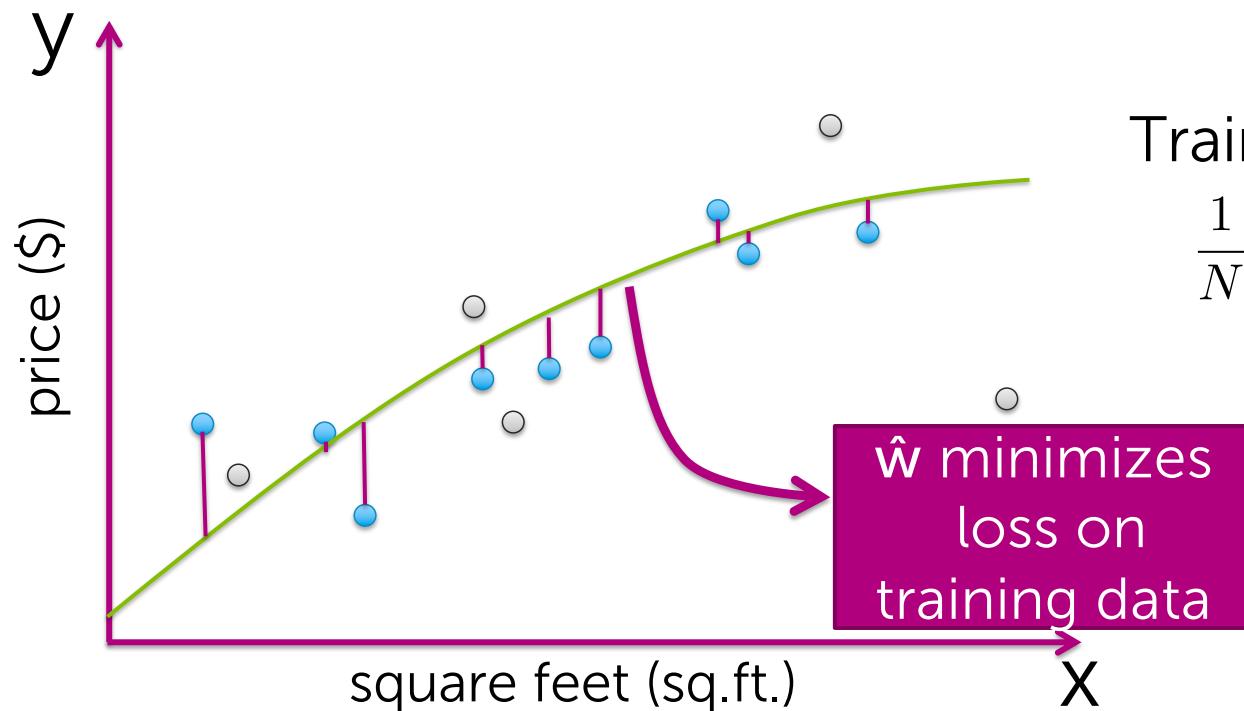
Define training data



Define training data



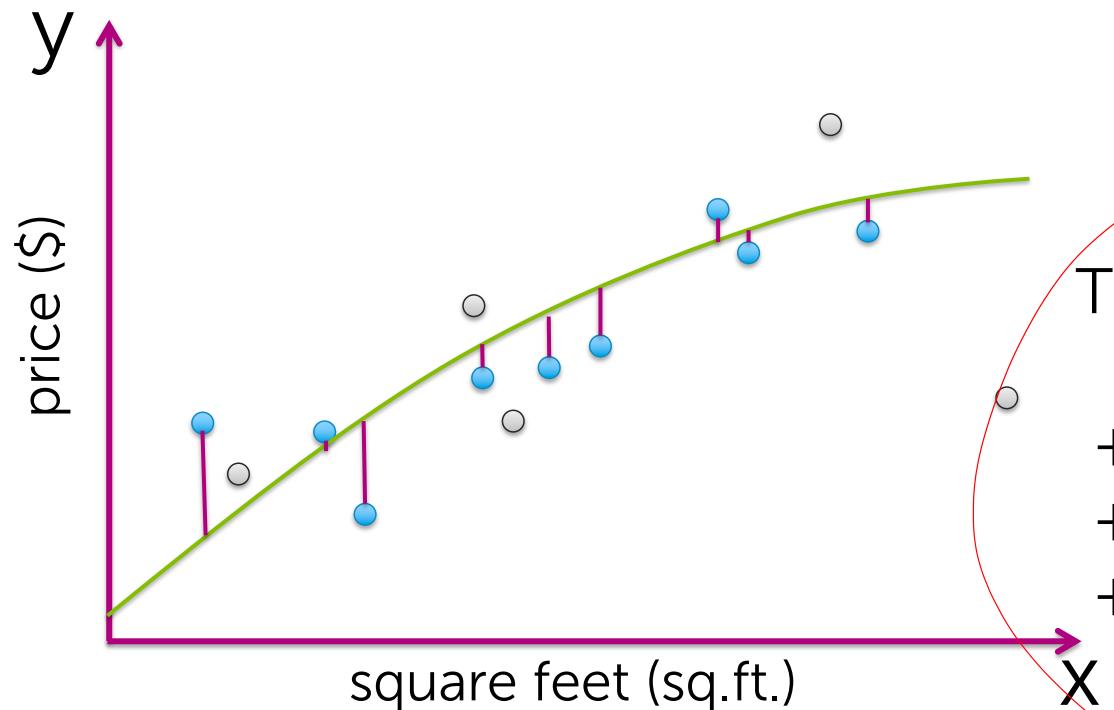
Example: Fit quadratic to minimize RSS



Training error (\hat{w}) =
$$\frac{1}{N} \sum_{i=1}^N (y_i - f_{\hat{w}}(x_i))^2$$

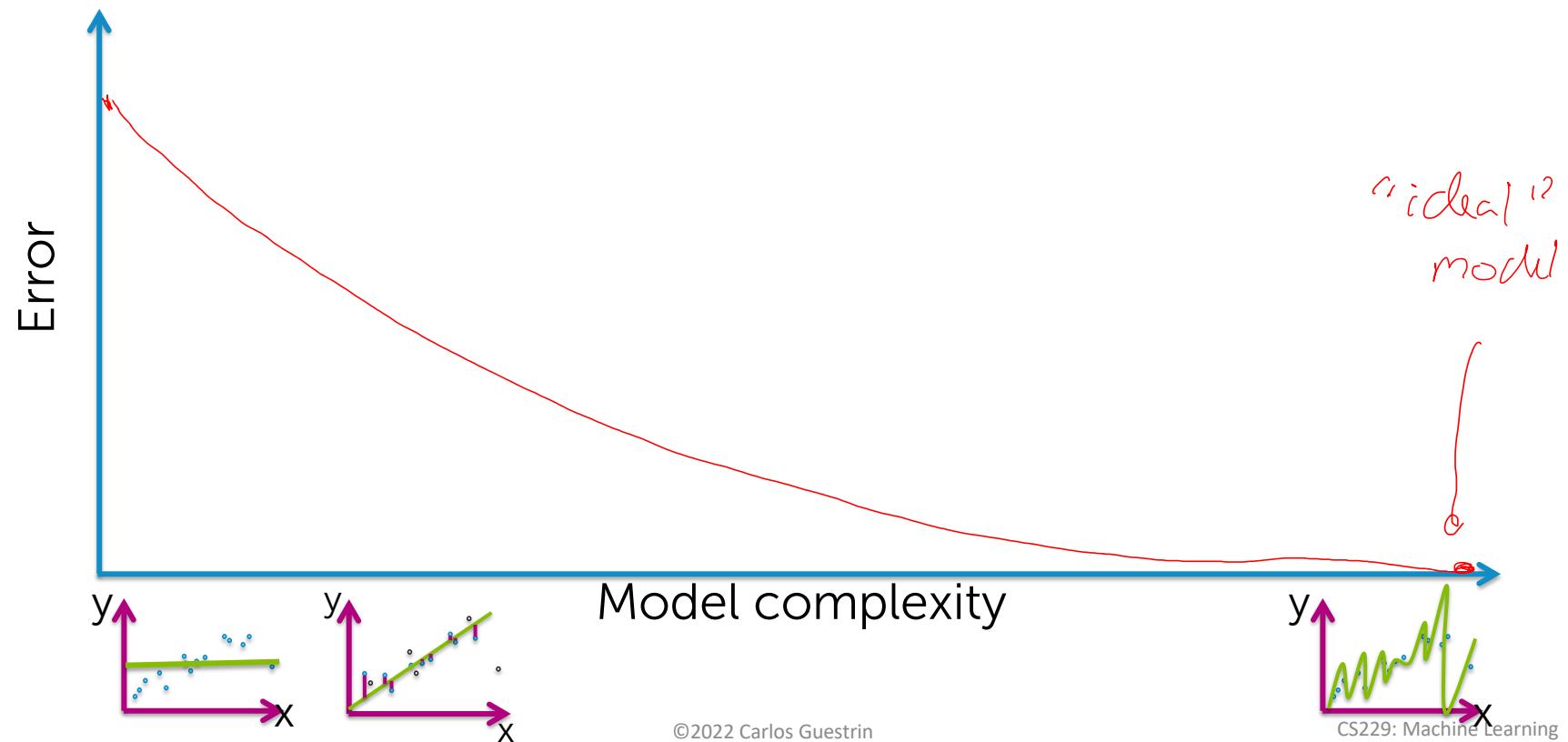
Example:

Use squared error loss $(y - f_{\hat{w}}(x))^2$



Training error (\hat{w}) = $1/N * [(\$_{train\ 1} - f_{\hat{w}}(\text{sq.ft.}_{train\ 1}))^2 + (\$_{train\ 2} - f_{\hat{w}}(\text{sq.ft.}_{train\ 2}))^2 + (\$_{train\ 3} - f_{\hat{w}}(\text{sq.ft.}_{train\ 3}))^2 + \dots \text{ include all training houses}]$

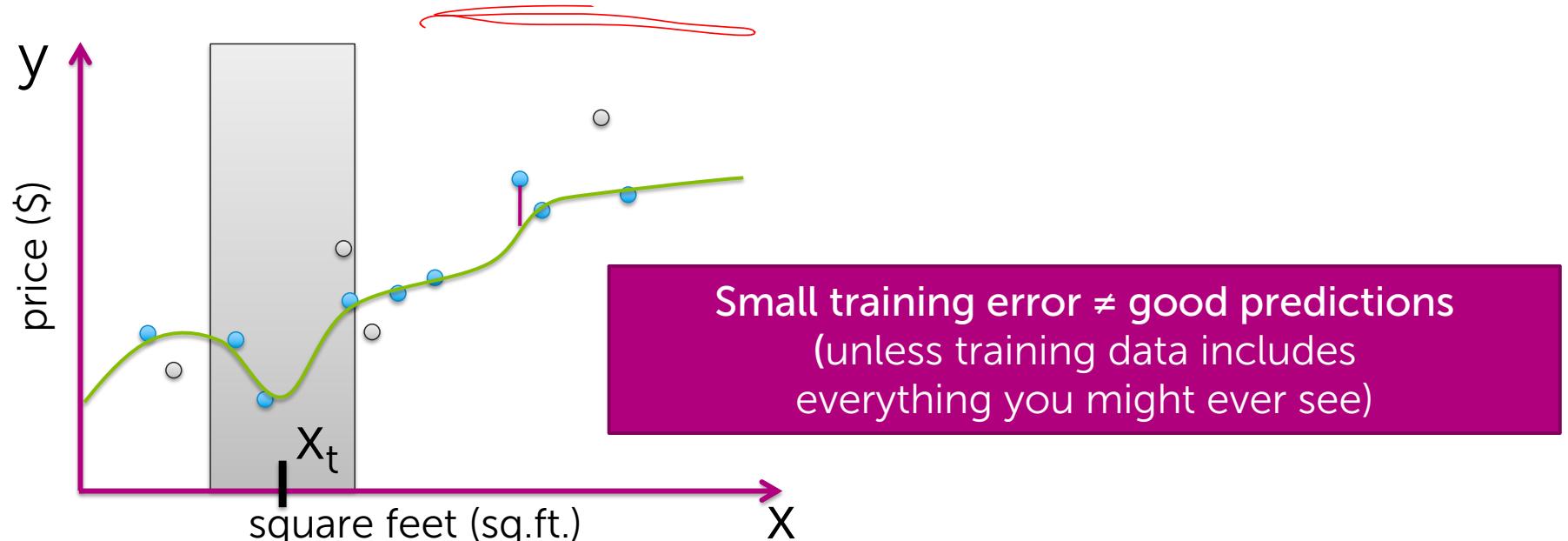
Training error vs. model complexity



Is training error a good measure of predictive performance?

Issue:

Training error is overly optimistic... \hat{w} was fit to training data

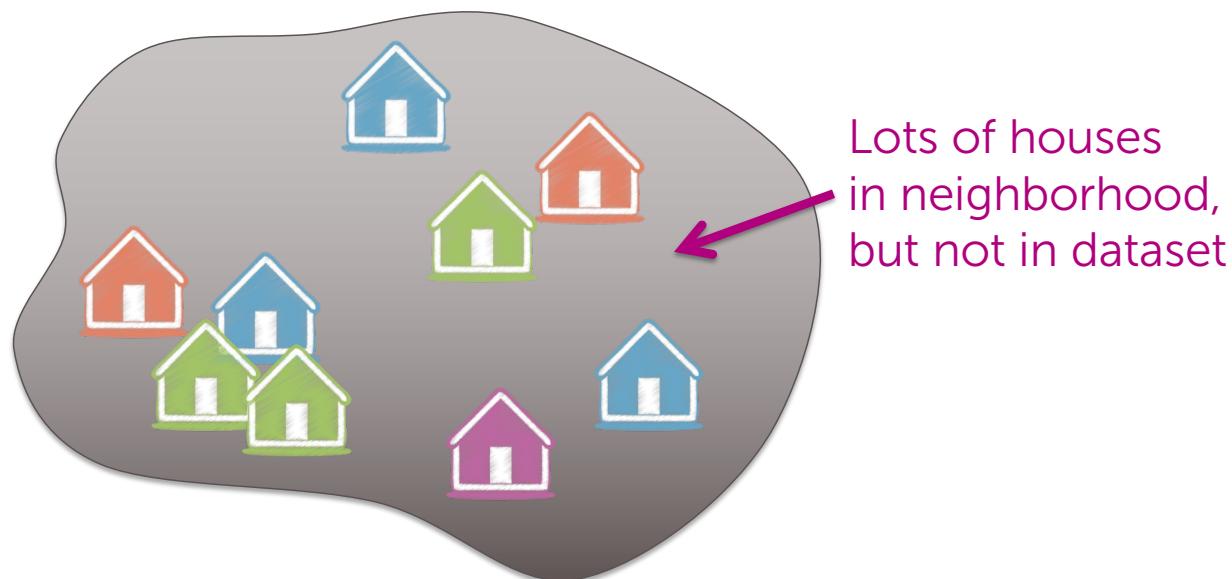


Assessing the loss

Part 2: Generalization (true) error

Generalization error

Really want estimate of loss over all possible (, ) pairs



Generalization error definition

Really want estimate of loss over all possible (, \$) pairs

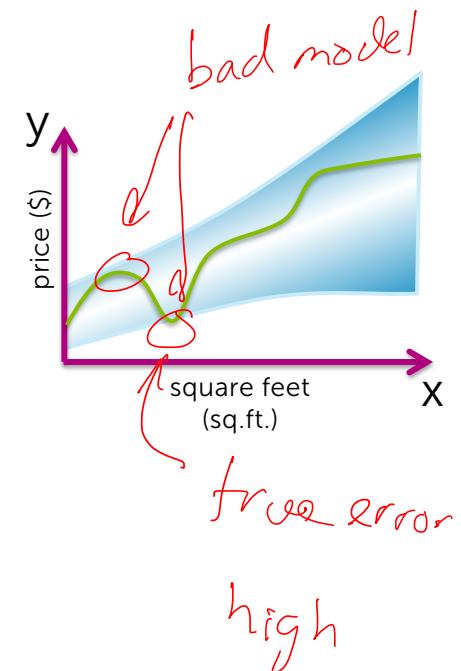
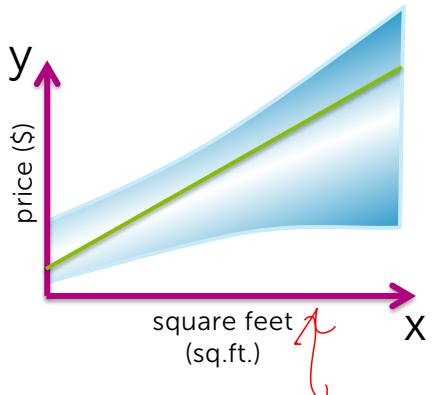
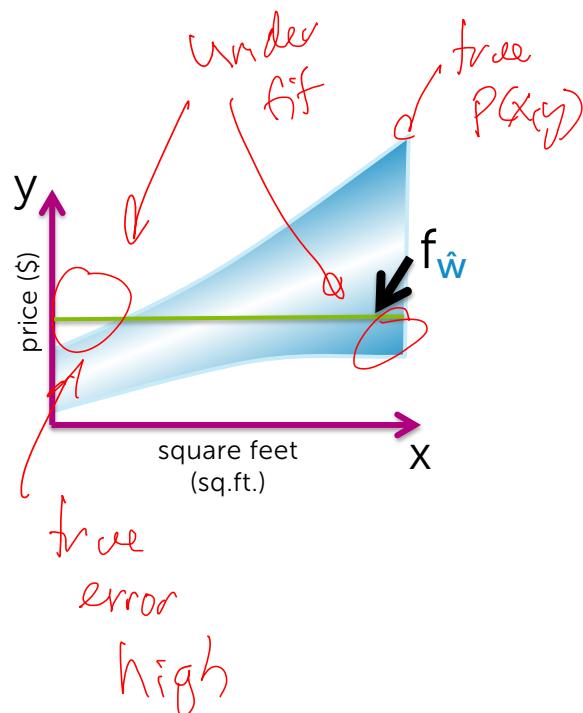
Formally:

average over all possible
(x, y) pairs weighted by
how likely each is

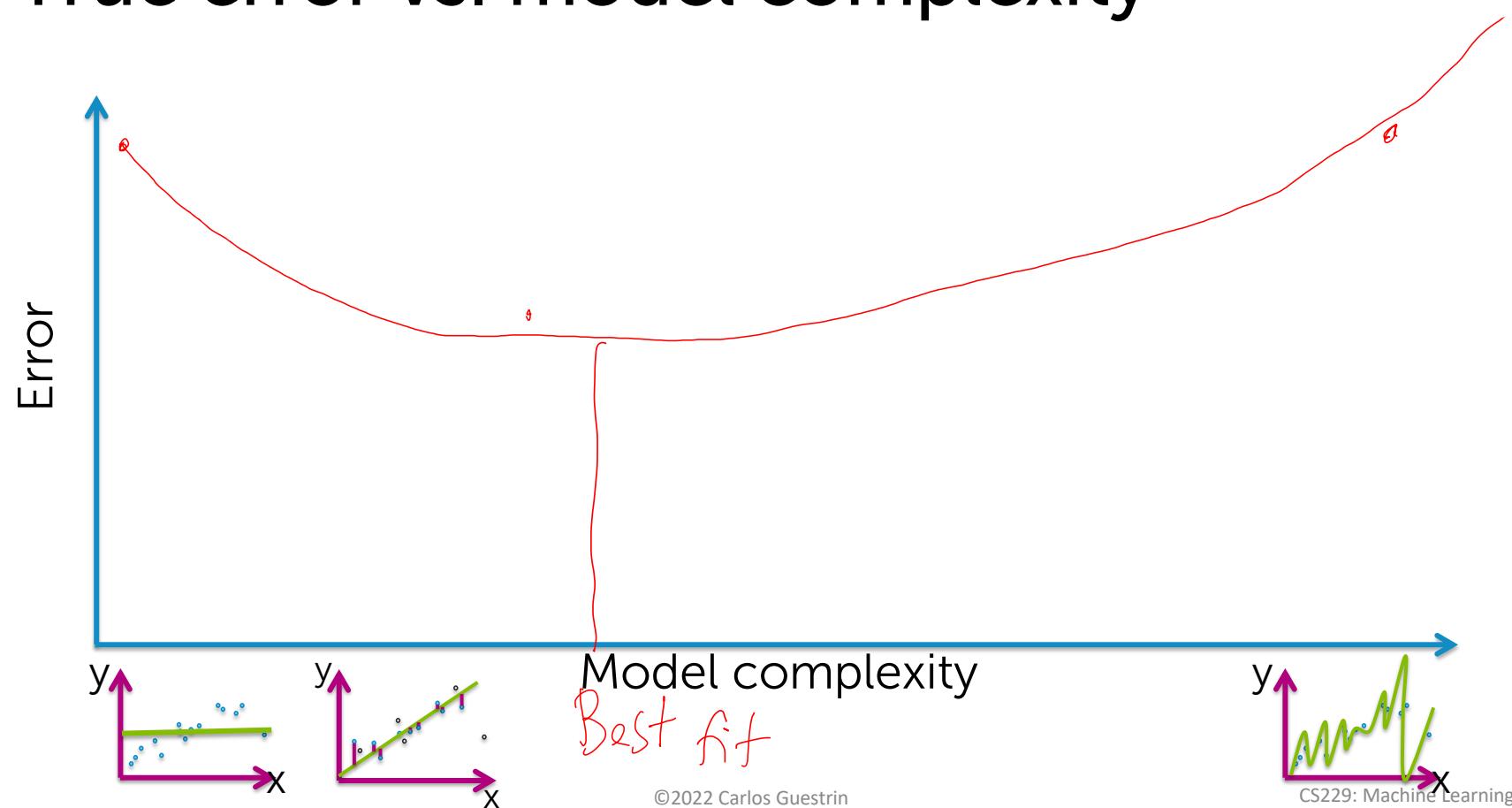
$$\text{generalization error} = E_{x,y} [L(y, f_{\hat{w}}(x))]$$

 
   
true prob. $p(x_i | y)$ fit using training data

Generalization error vs. model complexity



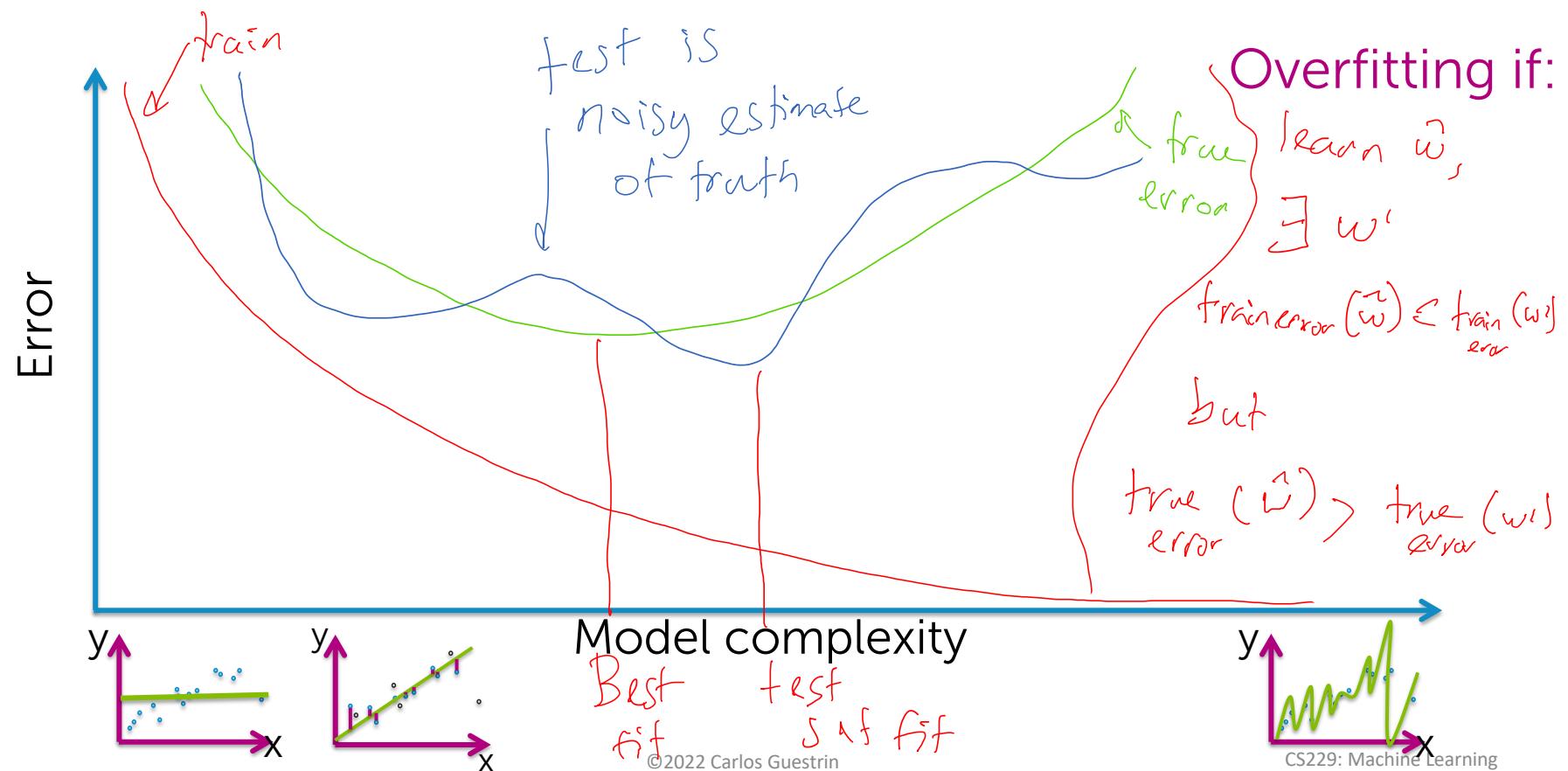
True error vs. model complexity



Assessing the loss

Part 3: Test error

Training, true, test error vs. model complexity



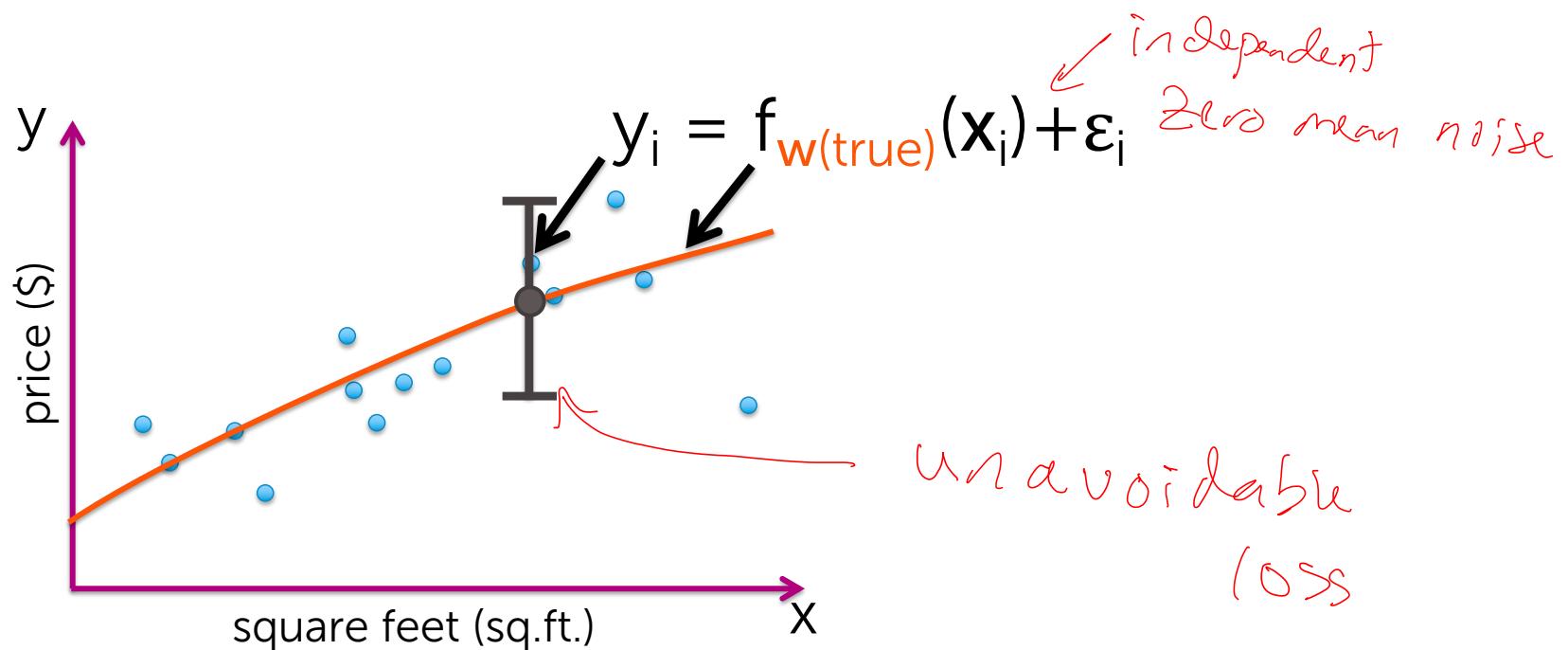
3 sources of error + the bias-variance tradeoff

3 sources of error

In forming predictions, there are 3 sources of error:

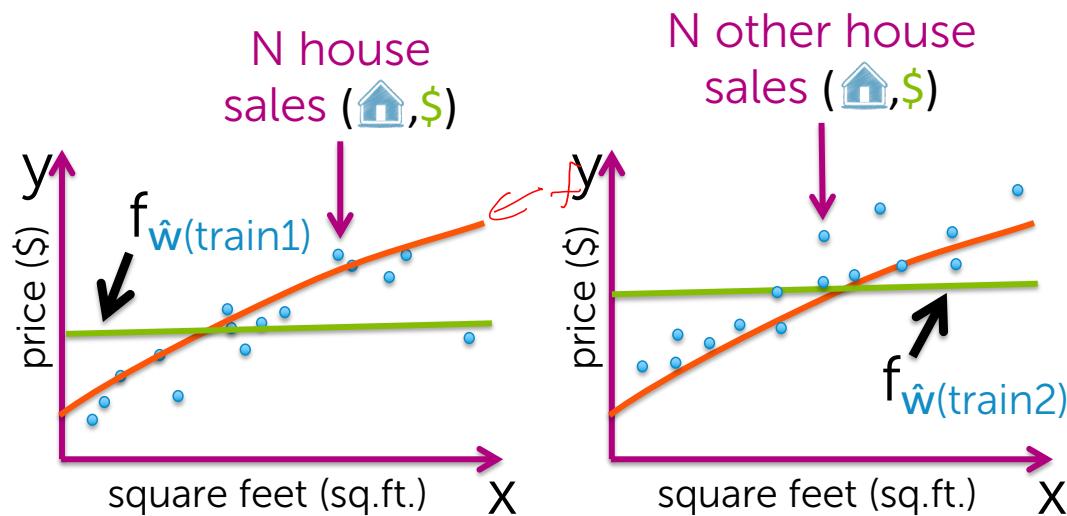
1. Noise
2. Bias ← how well can model fit data on avg.
3. Variance ← how much model would change if
 | change dataset samples

Data inherently noisy



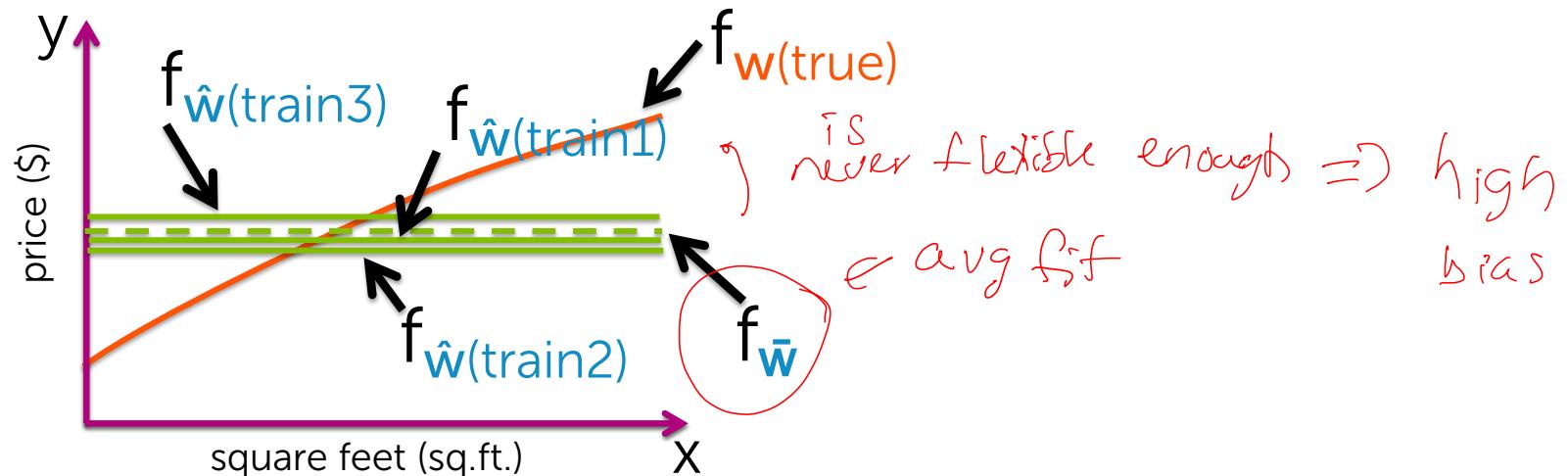
Bias contribution

Suppose we fit a constant function



Bias contribution

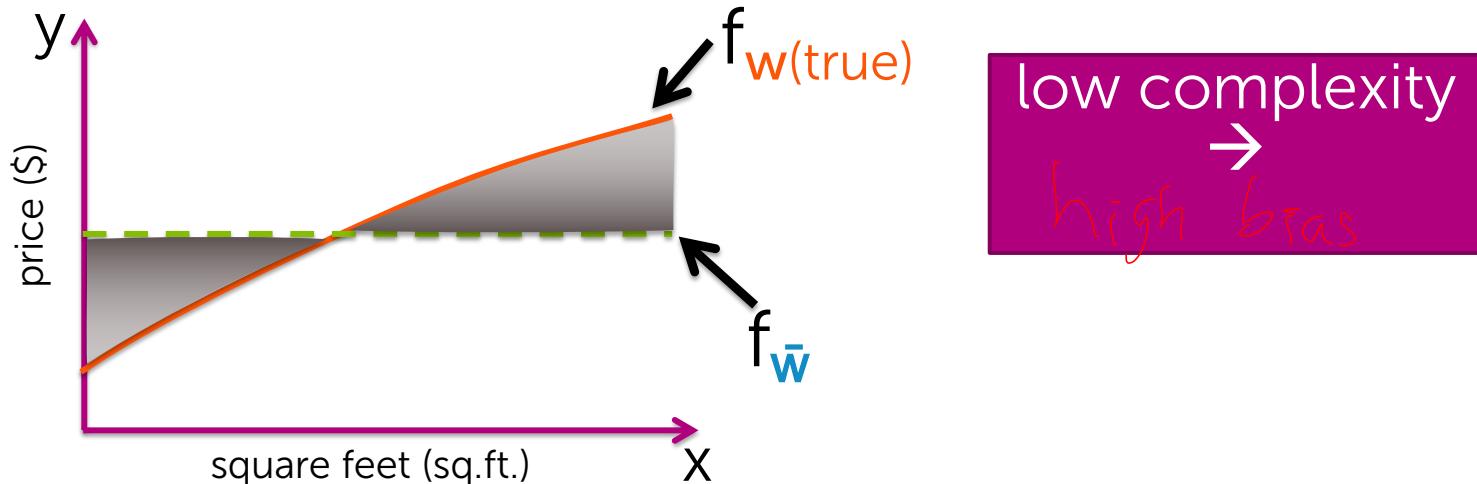
Over all possible size N training sets,
what do I expect my fit to be?



Bias contribution

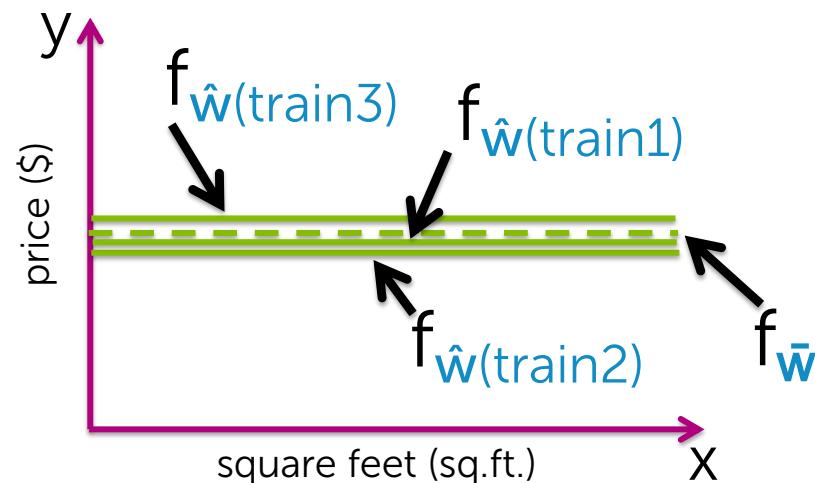
$$\text{Bias}(x) = f_{w(\text{true})}(x) - f_{\bar{w}}(x)$$

Is our approach flexible
enough to capture $f_{w(\text{true})}$?
If not, error in predictions.



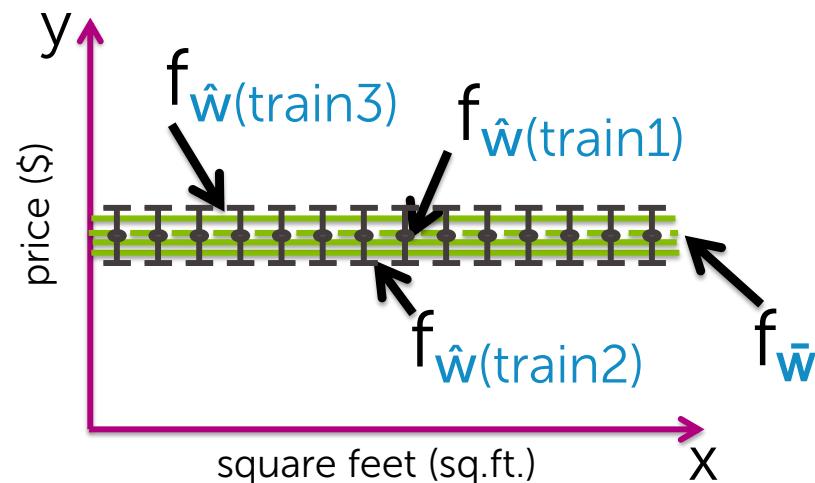
Variance contribution

How much do specific fits vary from the expected fit?



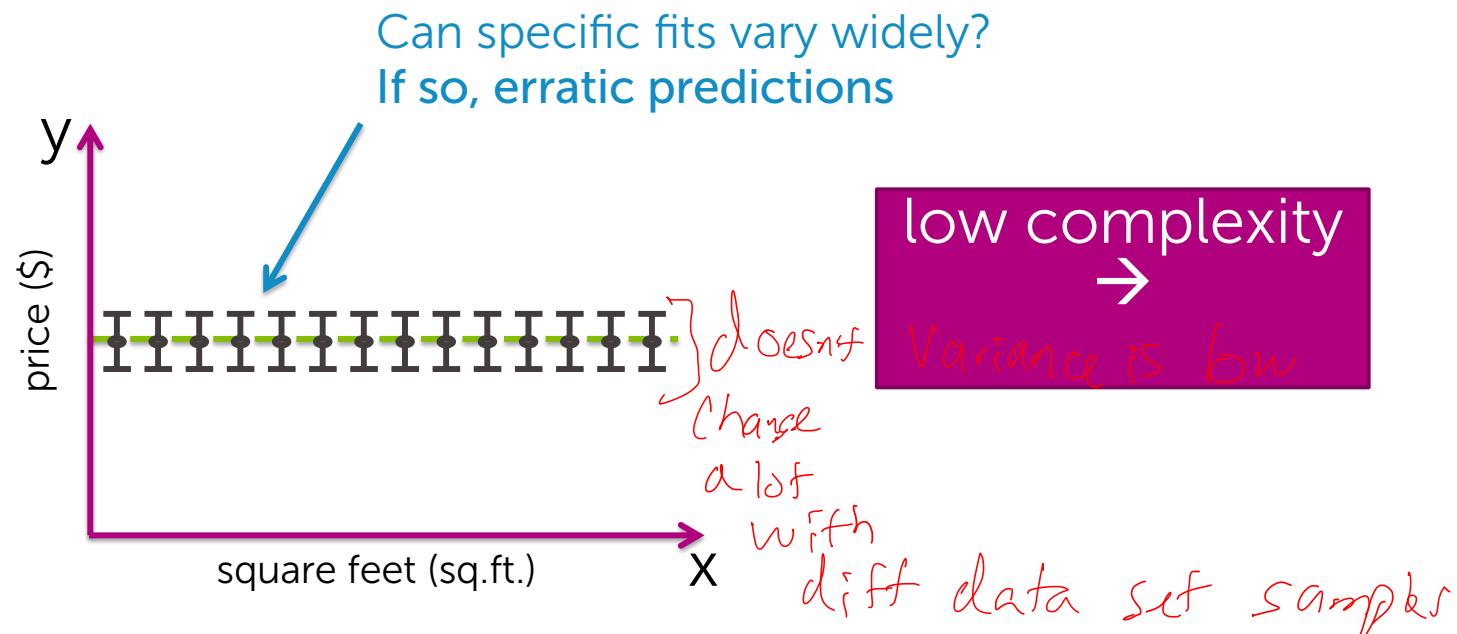
Variance contribution

How much do specific fits vary from the expected fit?



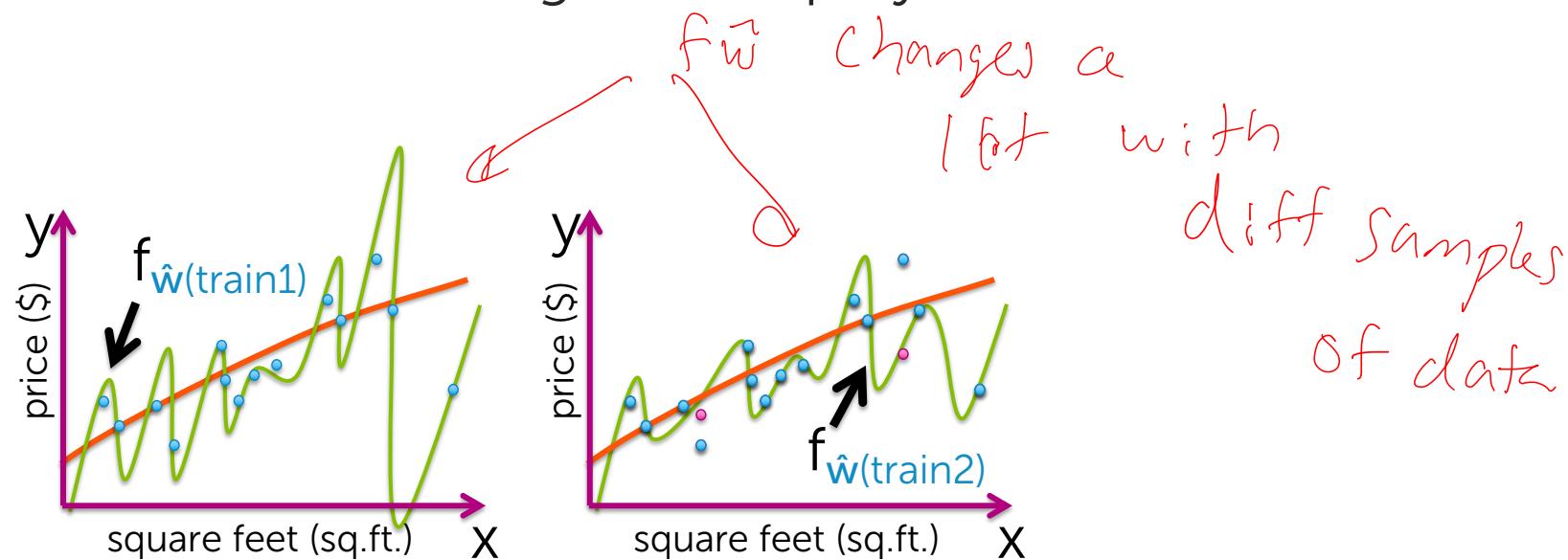
Variance contribution

How much do specific fits vary from the expected fit?



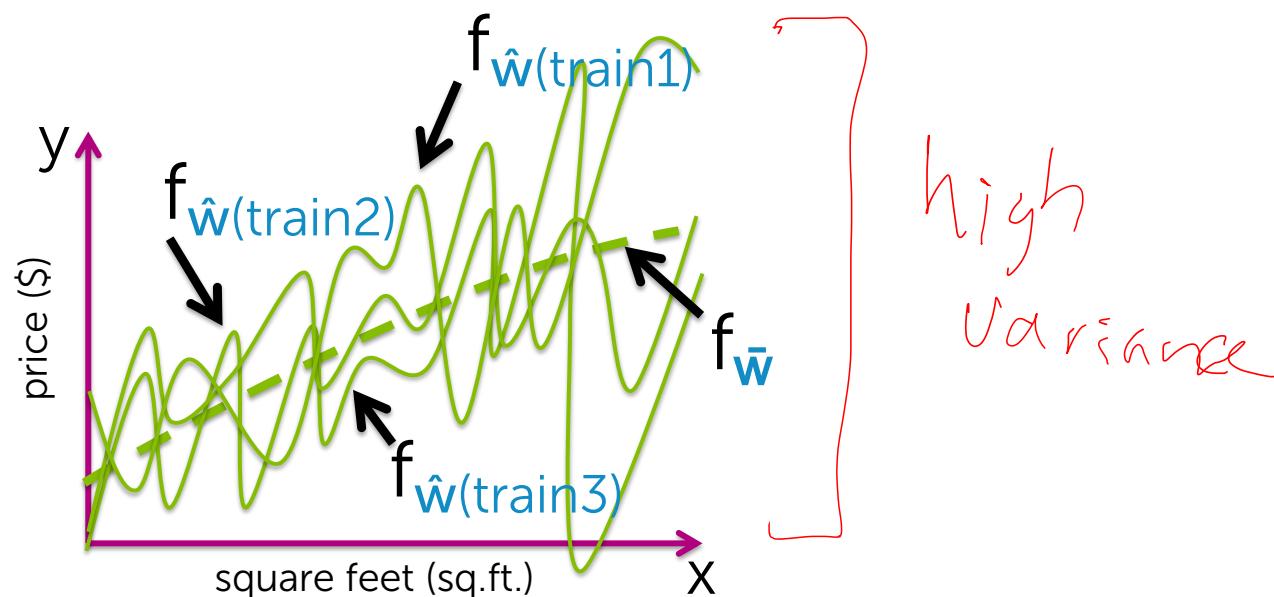
Variance of high-complexity models

Assume we fit a high-order polynomial

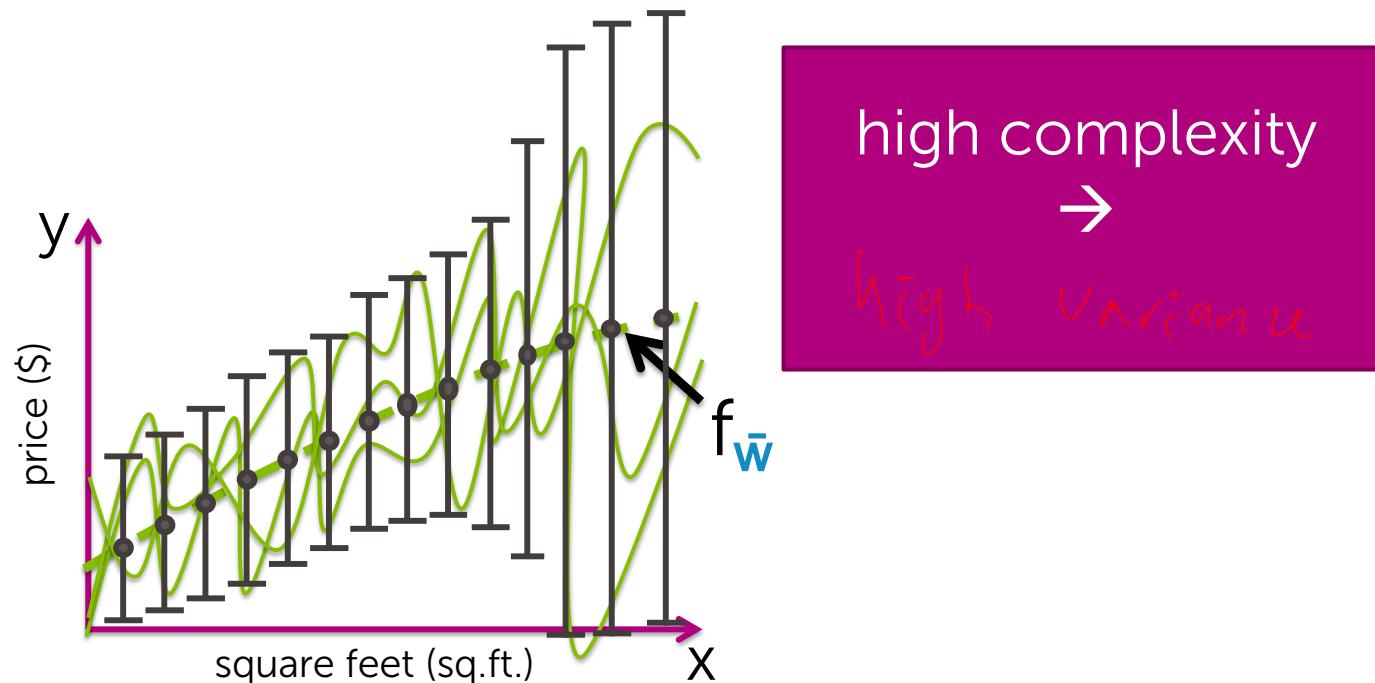


Variance of high-complexity models

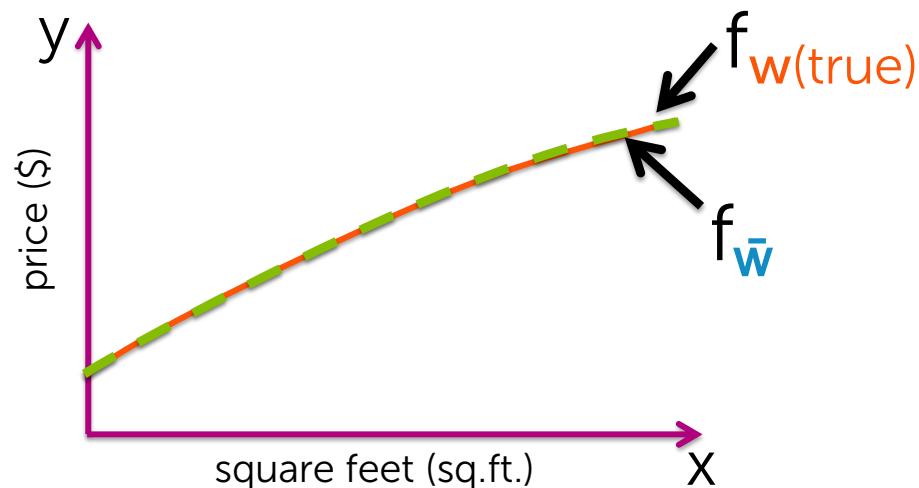
Suppose we fit a high-order polynomial



Variance of high-complexity models



Bias of high-complexity models

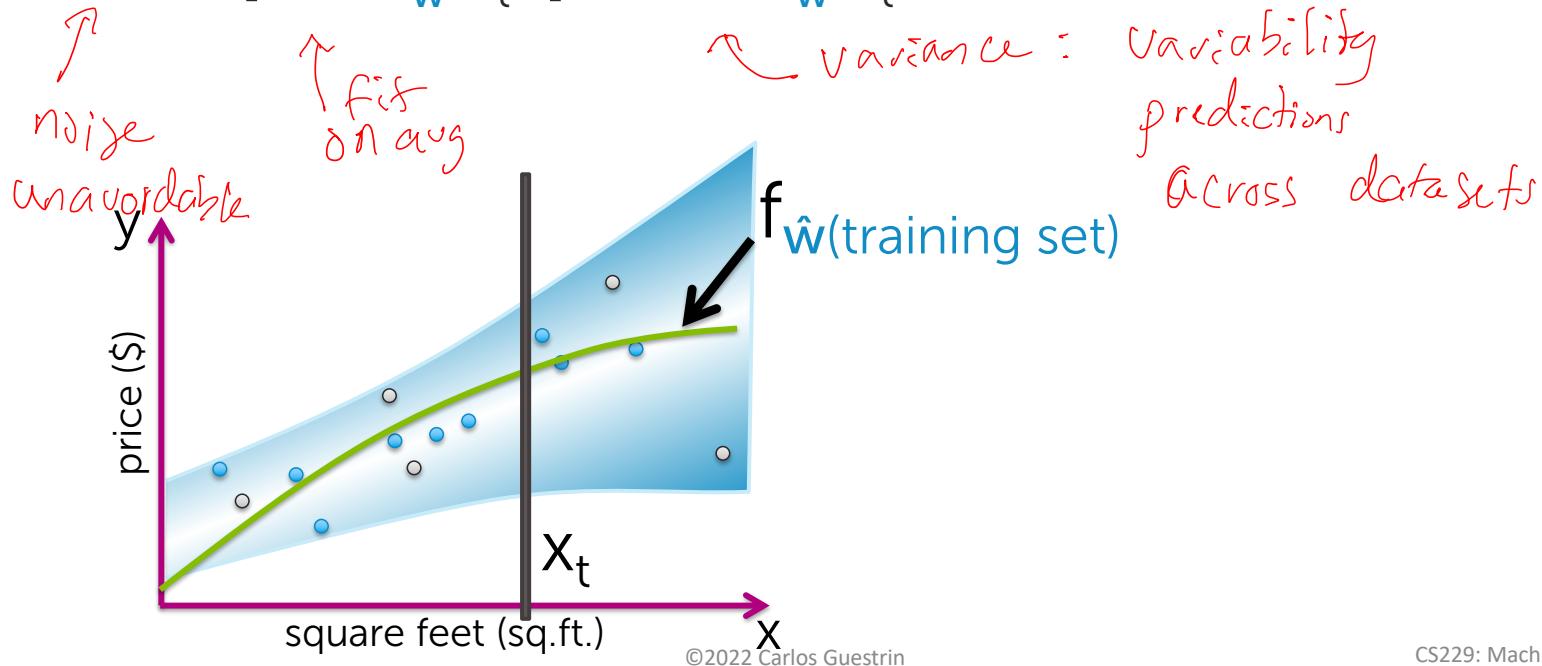


high complexity
→
low bias

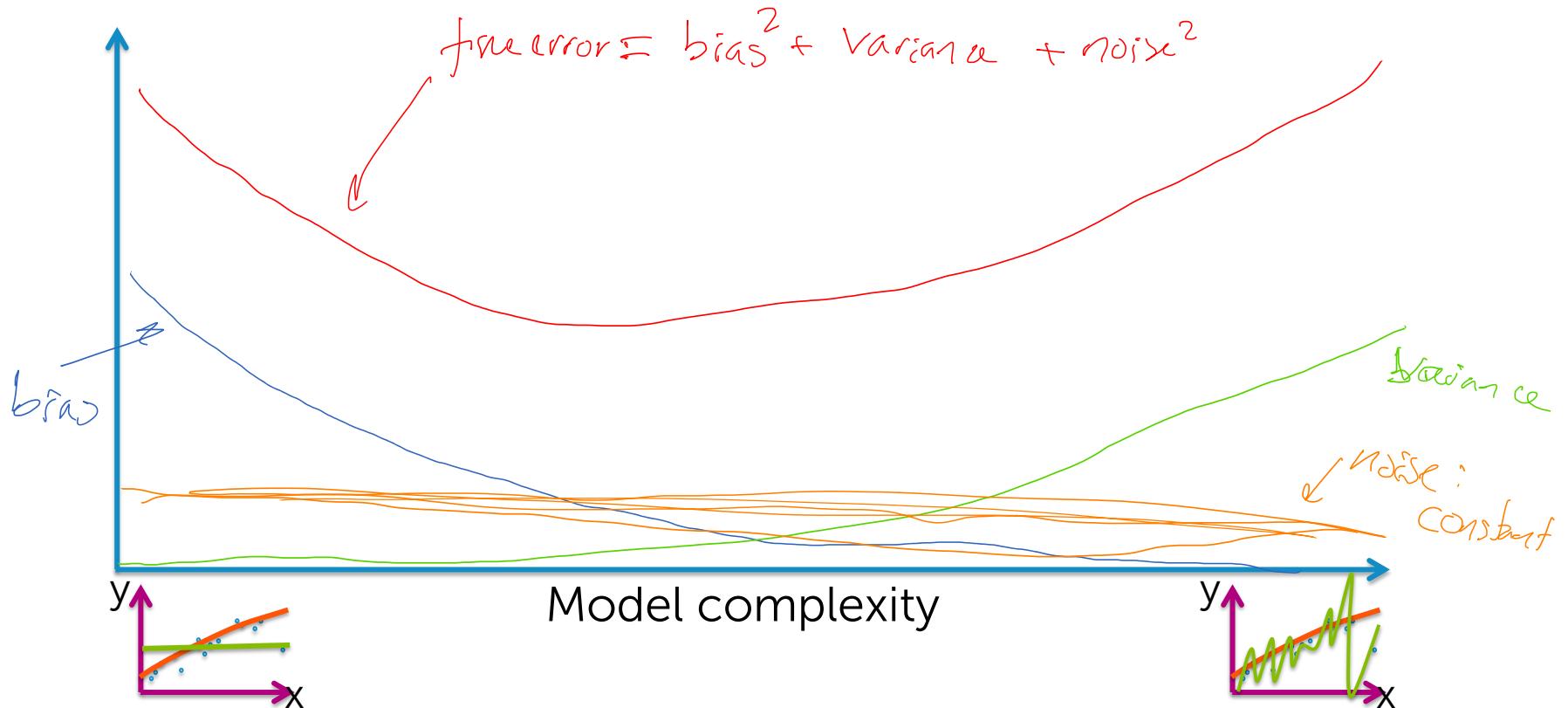
Sum of 3 sources of error

True errors
Average squared error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



Bias-variance tradeoff



Error vs. amount of data

for a fixed model complexity.
e.g., polynomials degree 73



Why 3 sources of error? A formal derivation

Deriving expected prediction error

Expected prediction error

$$\begin{aligned} &= E_{\text{train}} [\text{generalization error of } \hat{\mathbf{w}}(\text{train})] \\ &= E_{\text{train}} [E_{x,y} [L(y, f_{\hat{\mathbf{w}}(\text{train})}(x))]] \end{aligned}$$

1. Look at specific \underline{x}_t
2. Consider $L(y, f_{\hat{\mathbf{w}}}(x)) = \underline{(y - f_{\hat{\mathbf{w}}}(x))^2}$

Expected prediction error at x_t

$$= E_{\text{train}, y_t} [(y_t - f_{\hat{\mathbf{w}}(\text{train})}(x_t))^2]$$

Simplifying Notation

- Expected prediction error at x_t

$$= E_{\text{train}, y_t} [(y_t - f_{\hat{w}(\text{train})}(x_t))^2]$$

- Simple (and abusive 😊) notation:

– $y_t \rightarrow y$

– $f_{w(\text{true})}(x_t) \rightarrow f \leftarrow \text{true}$

– $f_{\hat{w}(\text{train})}(x_t) \rightarrow \hat{f} \leftarrow \text{learned}$

– $E_{\text{train}}[f_{\hat{w}(\text{train})}(x_t)] = \bar{f}_{\bar{w}}(x_t) \rightarrow \bar{f} \leftarrow \text{learned on avg.}$

Deriving expected prediction error

$$E[a+b] = E[a] + E[b]$$

a, b independent

Expected prediction error at x_t

$$= E_{\text{train}, y_t} [(y_t - f_{\hat{w}(\text{train})}(x_t))^2] = \underbrace{E_{\text{train}} [(y - \hat{f})^2]}_{} =$$

$$= E_{\text{train}} [(y - f) + (f - \hat{f})]^2$$

$$\begin{aligned} &= E[(y-f)^2] + 2 E_{\text{train}} [(y-f)(f-\hat{f})] + E_{\text{train}} [(f-\hat{f})^2] \\ &\quad \underbrace{\qquad}_{\text{by definition: noise}} \quad \underbrace{\qquad}_{\substack{0: \text{zero mean noise} \\ E_{\text{train}}[(y-f)] E_{\text{train}}[(f-\hat{f})]}} \quad \begin{array}{l} \xrightarrow{\triangle} \\ \text{MSE}(\hat{f}) \end{array} \\ &\quad \text{mean squared error} \\ &\approx \sigma^2 + \text{MSE}(\hat{f}) \end{aligned}$$

Equating MSE with bias and variance

$$\begin{aligned} \text{MSE}[f_{\hat{w}(\text{train})}(x_t)] &= E_{\text{train}}[(f - \hat{f})^2] \\ &= E_{\text{train}}[((f - \bar{f}) + (\bar{f} - \hat{f}))^2] \\ &= E_{\text{train}}[(f - \bar{f})^2] + 2E_{\text{train}}[(f - \bar{f})(\bar{f} - \hat{f})] + E_{\text{train}}[(\bar{f} - \hat{f})^2] \\ &\stackrel{\text{by defn.}}{=} \text{bias}^2 + 2E_{\text{train}}[(f - \bar{f})] \underbrace{E_{\text{train}}[\bar{f} - \hat{f}]}_{\bar{f}} - \underbrace{E[\bar{f}] - E[\hat{f}]}_{\delta} + \text{Var}(\hat{f}) \\ &= \text{bias}^2 + \text{Var} \end{aligned}$$

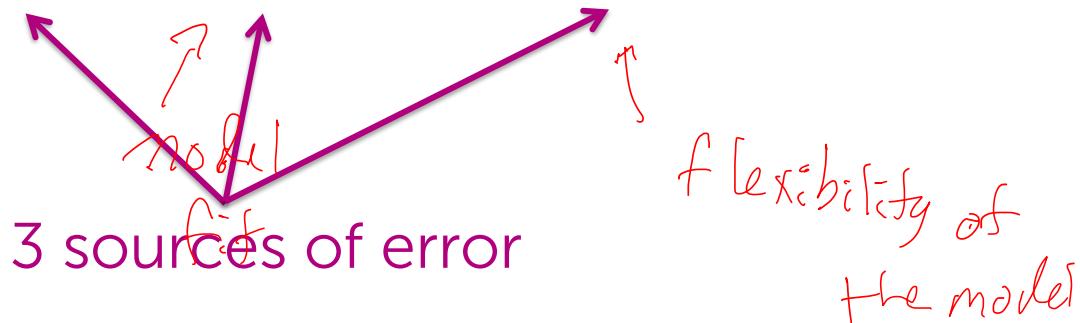
mean ↓ random ↓ var.

Putting it all together

Expected prediction error at x_t

$$= \sigma^2 + \text{MSE}[f_{\hat{w}}(x_t)]$$

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



Summary of bias-variance tradeoff

What you can do now...

- Contrast relationship between model complexity and train, true and test loss
- Compute training and test error given a loss function for different model complexities
- List and interpret the 3 sources of avg. prediction error
 - Irreducible error, bias, and variance